

---

# Embodied Agent Interface: Benchmarking LLMs for Embodied Decision Making

---

Manling Li<sup>1\*</sup>, Shiyu Zhao<sup>1\*</sup>, Qineng Wang<sup>1\*</sup>, Kangrui Wang<sup>1\*</sup>, Yu Zhou<sup>1\*</sup>,  
Sanjana Srivastava<sup>1</sup>, Cem Gokmen<sup>1</sup>, Tony Lee<sup>1</sup>, Li Erran Li<sup>2</sup>, Ruohan Zhang<sup>1</sup>, Weiyu Liu<sup>1</sup>,  
Percy Liang<sup>1</sup>, Li Fei-Fei<sup>1</sup>, Jiayuan Mao<sup>3</sup>, Jiajun Wu<sup>1</sup>

<sup>1</sup>Stanford University <sup>2</sup>Amazon <sup>3</sup>MIT

<https://cs.stanford.edu/~manlingl/projects/embodied-eval>

## Abstract

We aim to evaluate Large Language Models (LLMs) for embodied decision making. While a significant body of work has been leveraging LLMs for decision making in embodied environments, we still lack a systematic understanding of their performance because they are usually applied in different domains, for different purposes, and built based on different inputs and outputs. Furthermore, existing evaluations tend to rely solely on a final success rate, making it difficult to pinpoint what ability is missing in LLMs and where the problem lies, which in turn blocks embodied agents from leveraging LLMs effectively and selectively. To address these limitations, we propose a generalized interface (EMBODIED AGENT INTERFACE) that supports the formalization of various types of tasks and input-output specifications of LLM-based modules. Specifically, it allows us to unify 1) a broad set of embodied decision-making tasks involving both state and temporally extended goals, 2) four commonly-used LLM-based modules for decision making: goal interpretation, subgoal decomposition, action sequencing, and transition modeling, and 3) a collection of fine-grained metrics which break down evaluation into various types of errors, such as hallucination errors, affordance errors, various types of planning errors, etc. Overall, our benchmark offers a comprehensive assessment of LLMs' performance for different subtasks, pinpointing the strengths and weaknesses in LLM-powered embodied AI systems and providing insights for effective and selective use of LLMs in embodied decision making.

**1 Introduction**

Large Language Models (LLMs) have emerged as powerful tools for building embodied decision-making agents that can take natural language instructions from humans (such as “*cleaning the refrigerator*”, “*polishing furniture*”) and achieve the specified goals through a sequence of actions in various digital and physical environments. Despite many reports of their success, our understanding of LLMs’ full capabilities and limitations in embodied decision making remains limited. Existing evaluation methods fall short of providing a comprehensive insight due to three key limitations: the lack of standardization of 1) embodied decision-making tasks, 2) modules that an LLM can interface with or be implemented for, and 3) fine-grained evaluation metrics beyond a single success rate. In this paper, we propose EMBODIED AGENT INTERFACE, to address these challenges.

(1) **Standardization of goal specifications:** We want embodied agents to achieve goals. However, the specification of goals and the criteria for agents’ success evaluation vary significantly across

---

\*Equal contribution.

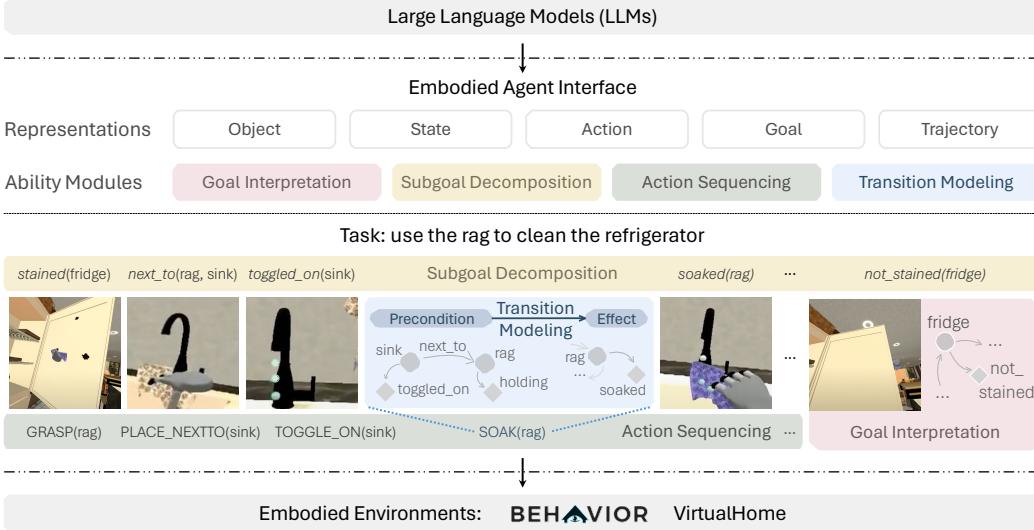


Figure 1: EMBODIED AGENT INTERFACE unifies a broad set of tasks involving both state and temporally extended goals and four LLM-based modules for decision making.

33 different domains, even for similar tasks. For example, BEHAVIOR [1] focuses on achieving a state  
 34 that satisfies certain *state goals* (e.g., “not *stained(fridge)*” in Figure 1), while VirtualHome [2] uses  
 35 temporally extended goals by imposing temporal order constraints on actions. We include an extended  
 36 discussion in Appendix N.1. Our EMBODIED AGENT INTERFACE implements a general object-  
 37 centric state and action representation, where object states, relations, and actions are represented in  
 38 abstract language terms (see Figure 1). Our innovation is to describe goals as linear temporal logic  
 39 (LTL) formulas, which define task-success criteria over trajectories. LTL affords the specification of  
 40 both state-based and temporally extended goals and allows for alternative goal interpretations.

41 (2) **Standardization of modules and interfaces:** Existing LLM-based embodied agent frameworks  
 42 often make different assumptions based on the availability of additional knowledge and external  
 43 modules. For instance, Code as Policies [3] and SayCan [4] utilize LLMs for action sequencing given  
 44 a given set of primitive skills, while LLM+P [5] uses LLMs for goal interpretation and generates  
 45 plans using PDDL planners with given domain definitions; Ada [6] leverages LLMs to generate  
 46 high-level planning domain definitions in PDDL and uses a low-level planner to generate control  
 47 commands. Consequently, they have defined different input-output specifications for the LLM  
 48 module, making comparisons and evaluations challenging. In EMBODIED AGENT INTERFACE,  
 49 built on top of our object-centric and LTL-based task specification, we formalize four critical *ability*  
 50 *modules* in LLM-based embodied decision making, as illustrated in Figure 1: *Goal Interpretation*,  
 51 *Subgoal Decomposition*, *Action Sequencing*, *Transition Modeling*. We formalize the input-output  
 52 specifications that LLMs can use to interface with other modules in the environment. This modular  
 53 interface automatically enables the integration of different LLM-based and external modules.

54 (3) **Broad coverage of evaluation and fine-grained metrics:** Current evaluations of LLMs for  
 55 embodied decision making have been overly simplified, usually focusing on the success rate of a  
 56 single task. The recent work LOTA-Bench [7] aims to break down the evaluation but is limited  
 57 to generating action sequences and does not support analysis of fine-grained planning errors. Our  
 58 EMBODIED AGENT INTERFACE, leveraging object-centric and factorized representations of states  
 59 and actions, implements a collection of fine-grained evaluation metrics, designed to automatically  
 60 locate different types of errors such as hallucination errors, different types of planning errors (e.g.,  
 61 object affordance errors, wrong action orders, etc.). Figure 2 illustrates different types of errors made  
 62 by GPT-4o on four different ability modules across two simulators.

63 We implement EMBODIED AGENT INTERFACE on two embodied decision-making benchmarks:  
 64 BEHAVIOR [1] and VirtualHome [2], and evaluated 15 different LLMs. Figure 2 visualizes the  
 65 performance of 5 representative LLMs on different tasks in Behavior. Our key findings are

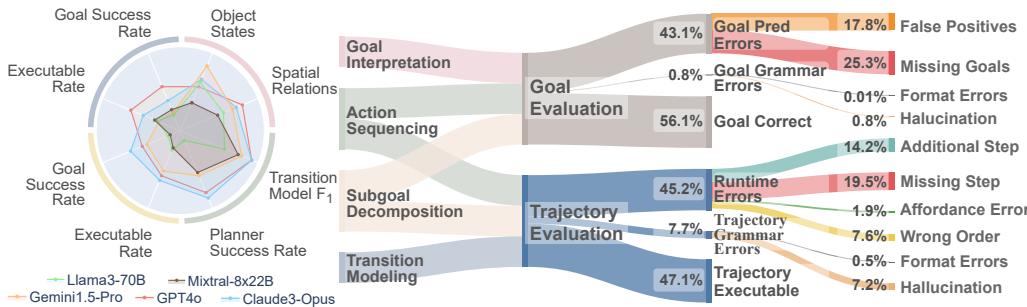


Figure 2: EMBODIED AGENT INTERFACE supports a collection of fine-grained metrics and provides automatic toolkits for error analysis and benchmarking different LLMs on various embodied decision-making tasks.

- Most LLMs struggle to faithfully translate natural language instructions into grounded states (objects, object states, and relations) in the environment. They sometimes predict intermediate subgoals as part of the final goals, e.g., predicting the state *open*(freezer) for task “*drinking water*”.
- Reasoning ability is a crucial aspect that LLMs should improve. Trajectory feasibility errors are common (45.2%), with a large portion of missing step (19.5%) and additional step (14.2%) errors, often due to overlooking preconditions. For instance, LLMs may ignore the agent’s *sitting* or *lying* state and fail to include a *standup* action before executing other actions. They sometimes also fail to understand the need to *open* a *closed* object before *fetching* items from inside. Additional step errors frequently occur when LLMs output actions for previously achieved goals.
- Trajectory evaluation performance decreases as the trajectory sequence length increases; goal evaluation performance decreases when the environment becomes more complex, involving a larger variety of object and state features.
- LLM errors include not only hallucinations of nonexistent objects and actions but also a heavy reporting bias. They often focus on surface forms but ignore details. For example, “put the turkey on the table” should be interpreted as “put the turkey on a plate, and place the plate on the table.”
- Subgoal decomposition is not strictly easier than action sequencing in abstract action spaces.
- We further provide quantitative analysis for the robustness of the modules through sensitivity analysis, pipeline-based versus modularized comparison, and replanning. These analyses aim to identify potential ways to integrate LLM-based and external modules.

## 2 Embodied Agent Interface Based on LTL

Table 1 summarizes our EMBODIED AGENT INTERFACE. First, we define an **embodied decision-making problem representation**  $\langle \mathcal{U}, \mathcal{S}, \mathcal{A}, g, \phi, \bar{a} \rangle$ , which is a language-based, object-centric abstraction for embodied agent environments with *objects* ( $o \in \mathcal{U}$ ), *states* ( $s \in \mathcal{S}$ ), *actions* ( $a \in \mathcal{A}$ ), *goal*  $g$ , *subgoal*  $\phi$ , and trajectories  $\bar{a}$ . Second, we formally define four **ability modules**  $\langle \mathcal{G}, \Phi, \mathcal{Q}, \mathcal{T} \rangle$ , including their standardized input-output specifications. They are fundamental and commonly-used modules that LLMs can be implemented for and interface with: the *goal interpretation* module  $\mathcal{G}$ , the *action sequencing* module  $\mathcal{Q}$ , the *subgoal decomposition* module  $\Phi$ , and the *transition modeling* module  $\mathcal{T}$ . In this paper, we focus on object-centric modeling: states are described as relational features among entities in the environment, actions are defined functions that take entity names as inputs and can be executed in the environment, goals and subgoals are defined as linear-temporal logic (LTL) formulas on states and actions. We define each component in detail as follows.

### 2.1 Representation for Objects, States and Actions

In EMBODIED AGENT INTERFACE, a state is represented as a tuple  $s = \langle \mathcal{U}, \mathcal{F} \rangle$ , where  $\mathcal{U}$  is the universe of objects, assumed to be a fixed finite set.  $\mathcal{F}$  is a set of relational Boolean features. Each  $f \in \mathcal{F}$  is a table where each entry is associated with a tuple of objects  $(o_1, \dots, o_k)$ . Each entry has the value of the feature in the state, and  $k$  is the arity of the feature. Actions can be viewed as primitive functions that take objects as inputs, denoted as  $\langle name, args \rangle$ . Throughout the paper, we focus on tasks where states and actions are described in abstract language forms, including object states (e.g., *is-open*(cabinet1)), relations (e.g., *is-on*(rag0, window3)), and actions (e.g., *soak*(rag0)).

Table 1: Summary of notations used in EMBODIED AGENT INTERFACE.

	<b>Notation</b>	<b>Symbol</b>	<b>Description</b>
	Object	$u \in \mathcal{U}$	An object, which has relational features $f$
	State	$s = \langle \mathcal{U}, \mathcal{F} \rangle \in \mathcal{S}$	A tuple of the universe of objects and relational features
	Action	$a = \langle name, args \rangle \in \mathcal{A}$	A tuple of the action name and arguments
	Operator	$o = \langle name, vars \rangle \in \mathcal{O}$	An action schema: a tuple of the name and a list of parameters. Each $o$ can be instantiated into an action $a$
	Transition Model	$\mathcal{M} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$	The deterministic transition function of the environment
Environment Representations	Natural Language Goal	$g_{nl}$	A sentence in English
	LTL Goal	$g$	An LTL formula. Here, we only consider formulas containing a sequence of action items and a conjunction of propositions (for the final state): $g = a_1 \text{ then } \dots \text{ then } a_k \text{ then } (p_1 \wedge \dots \wedge p_\ell)$ .
	Action Trajectory	$\bar{a} = \{a_i\}_{i=1}^n$	A sequence of $n$ actions
	Subgoal Trajectory	$\bar{\phi} = \{\phi_i\}_{i=1}^m$	A sequence of LTL subgoals $\phi_i$ connected by “then”
	State-action Trajectory	$\bar{t} = \langle \{s_i\}_{i=0}^n, \{a_i\}_{i=1}^n \rangle$	A sequence of state-action pairs. $\forall t. s_{t+1} = \mathcal{M}(s_t, a_t)$
Abilities	Task	$\langle s_0, g, g_{nl} \rangle$	A tuple of the initial state and the LTL/Natural Language goals
	Goal Interpretation	$\mathcal{G} : \langle s_0, g_{nl} \rangle \rightarrow g$	Initial State & Natural Language Goal $\rightarrow$ LTL Goal
	Subgoal Decomposition	$\Phi : \langle s_0, g \rangle \rightarrow \bar{\phi}$	Initial State & Goal $\rightarrow$ Subgoal Trajectory
	Action Sequencing	$\mathcal{Q} : \langle s_0, g \rangle, \mathcal{M} \rightarrow \bar{a}$	Initial State & Goal & Transition Model $\rightarrow$ Action Trajectory
	Transition Modeling	$\mathcal{T} : \langle s_0, g \rangle, o \rightarrow \langle pre, eff \rangle$	Initial State & Goal & Operator $\rightarrow$ Preconditions & Effects

## 105 2.2 Representation for Goals, Subgoals, Action Sequences, and State-Action Trajectories

106 In EMBODIED AGENT INTERFACE, goals  $g$ , subgoals  $\phi$ , and action sequences  $\bar{a}$  are modeled as  
 107 linear temporal logic (LTL) formulas. This is motivated by two critical desiderata. First, we need an  
 108 expressive and compact language to describe both state-based and temporally extended goals. Second,  
 109 we need a unified interface between different LLM-based modules. LTL addresses both challenges.  
 110 At a high level, an LTL formula can describe state constraints (e.g., a subgoal should be achieved),  
 111 action constraints (e.g., a particular action should be executed), and possible temporal orders among  
 112 them (e.g., all dishes should be cleaned before we cook). By combining temporal connectives (such  
 113 as “eventually”) and propositional logic connectives (such as “or”), we can also flexibly describe  
 114 alternative goals or trajectories. As a byproduct, using a single description language for all inputs and  
 115 outputs enables us to design a unified metric to measure accuracy, which we detail in Appendix ??.

116 In EMBODIED AGENT INTERFACE, we use a fragment of the full linear temporal logic (LTL)  
 117 formalism on finite trajectories. We allow two types of atomic propositions: state propositions (object  
 118 properties and relations) and action propositions. Our LTL language contains Boolean conjunction  $\wedge$ ,  
 119 disjunction  $\vee$ , negation  $\neg$ , implication  $\Rightarrow$ , first-order logic quantifiers  $\forall, \exists, \exists^{=n}$  (the equal quantifier:  
 120 there are exactly  $n$  objects satisfying a condition), and the temporal connective **then**.

121 An LTL formula is a trajectory classifier semantically: the function  $eval(\phi, \bar{t})$  evaluates an LTL  
 122 formula  $\phi$  on a state-action sequence  $\bar{t}$ . We say that the state-action sequence satisfies  $\phi$  if  $eval(\phi, \bar{t}) =$   
 123 *true* (i.e., the goal  $\phi$  is satisfied). For state formulas  $\phi$  (formulas without **then**), we define  $eval(\phi, \bar{t}) =$   
 124  $\exists t. \phi(s_t)$  (“eventually” the goal is satisfied). For formulas connected by **then**,  $eval(\phi_1 \text{ then } \phi_2, \bar{t}) =$   
 125  $\exists k. \phi_1(\bar{t}_{\leq k}) \wedge \phi_2(\bar{t}_{>k})$  ( $\phi_2$  is achieved after  $\phi_1$ ), where  $\bar{t}_{\leq k}$  and  $\bar{t}_{>k}$  denote prefixes and suffixes.  
 126 Currently, we have not implemented other temporal connectives such as “globally” and “until” but  
 127 our overall framework can be extended to them. An LTL formula example of a subgoal plan for task  
 128 *browse Internet* is: “*ontop(character, chair) then holds\_rh(character, mouse)  $\wedge$  holds\_lh(character,*  
 129 *keyboard) then facing(character, computer)*”. We include LTL details in Appendix ??.

## 130 2.3 Ability Module 1: Goal Interpretation $\mathcal{G} : \langle s_0, g_{nl} \rangle \rightarrow g$

131 **Input-Output Specification.** The *goal interpretation* module takes the state  $s_0$  and a natural language  
 132 instruction  $g_{nl}$  as input, and generates an LTL goal  $\hat{g}$ , as a formal goal specification which a symbolic  
 133 planner can conceivably take as input. In this paper, we only generate simple LTL goals formed by  
 134 an ordered action sequence and a conjunction of propositions to be satisfied in the final state.

135 **Evaluation Metric.** An LTL goal can be evaluated by directly comparing it with the ground truth  
 136 goal  $g$ . While we have restricted generated  $\hat{g}$  to be simple LTL goals, we do not require the ground

137 truth goal  $g$  to be simple. Therefore, we additionally define  $\mathcal{G}$  that takes the object universe  $\mathcal{U}$  as  
 138 input to translate  $g$  to a set of simple LTL goals  $g_0, g_1, \dots, g_k$  where all  $g_i$ 's entail  $g$ . We describe  
 139 our implementation in the Appendix. Given two simple LTL goals  $g_i$  and  $\hat{g}$ , the accuracy of  $\hat{g}$  can be  
 140 computed as an  $F_1$  set-matching score between them. Let  $g = a_1 \text{ then } a_k \text{ then } (p_1 \wedge \dots \wedge p_\ell)$ . We  
 141 define  $\text{set}(g) = \{\{a_i\}_{i=1}^k\} \cup \{p_i\}_{i=1}^\ell$  (i.e., the action sequence  $\{a_i\}$  is treated as a single element).  
 142 The  $F_1$  score between  $g$  and  $\hat{g}$  is defined as:  $F_1(g, \hat{g}) = \max_{g_i \in \mathcal{G}(g, \mathcal{U})} F_1(\text{set}(g_i), \text{set}(\hat{g}))$ .

#### 143 2.4 Ability Module 2: Action Sequencing $\mathcal{Q} : \langle s_0, g \rangle, \mathcal{M} \rightarrow \bar{a}$

144 **Input-Output Specification.** The *action sequencing* module takes the task  $\langle s_0, g \rangle$  as input, and the  
 145 transition model  $\mathcal{M}$ , and generates an action sequence  $\bar{a} = \{a_i\}_{i=1}^n$ .

146 **Evaluation Metric.** We use two evaluation metrics for the action sequencing module. First, the  
 147 *trajectory feasibility evaluation* focuses on evaluating whether the trajectory is executable (i.e., all  
 148 actions are feasible). We will execute the trajectory  $\bar{a}$  from  $s_0$  in the simulator. When infeasible  
 149 action presents, the execution may stop at an early step and we categorize the execution failure into  
 150 missing steps, additional steps, wrong temporal order, and affordance errors.

151 Second, the *goal satisfaction evaluation* evaluates if the goal is satisfied after executing  $\bar{a}$ . Specifically,  
 152 we obtain  $T = \langle \{s_i\}_{i=0}^m, \{a_i\}_{i=1}^m \rangle$  by executing  $\bar{a}$ , and directly use the  $\text{eval}(g, T)$  function to check  
 153 for goal satisfaction. We also evaluate the *partial goal satisfaction evaluation*, which is the percentage  
 154 of “subgoals” in  $g$  that are satisfied in  $\bar{a}$ . To compute this partial success rate, we again consider  
 155 all simple LTL goals  $g_i$  derived from  $g$ . Let  $g_i = a_1 \text{ then } a_k \text{ then } (p_1 \wedge \dots \wedge p_\ell)$ . If there is a  
 156 subsequence in  $\bar{a}$  that is the same as  $\{a_j\}_{j=1}^k$ , we consider the action sequence successfully executed.  
 157 Next, we evaluate all final state propositions  $p_j$  and give models partial credits based on the number  
 158 of propositions satisfied in  $s_m$ . Finally,  $\text{PartialSucc}(\bar{a}, g) = \max_{g_i \in \mathcal{G}(g, \mathcal{U})} \text{PartialSucc}(\bar{a}, g_i)$ .

#### 159 2.5 Ability Module 3: Subgoal Decomposition $\Phi : \langle s_0, g \rangle \rightarrow \bar{\phi}$

160 **Input-Output Specification.** The *subgoal decomposition* module takes the task  $\langle s_0, g \rangle$  as input  
 161 and generates a sequence of subgoals  $\bar{\phi} = \{\phi_i\}_{i=1}^k$ , where each  $\phi_i$  is an LTL formula. The entire  
 162 sequence  $\bar{\phi}$  can also be represented as a single LTL formula.

163 **Evaluation Metric.** To evaluate the subgoal decomposition module, we use a customized planner  
 164 to refine it into an action sequence  $\bar{a}$ . This subgoal-action mapping function  $\mathcal{AM}(\bar{\phi}, s_0)$  takes the  
 165 LTL representation of  $\bar{\phi}$  and  $s_0$  and generates a state-action sequence  $\bar{t}$ . We implement this with a  
 166 breadth-first search. Then, we use the same metrics in *action sequencing* for evaluation: trajectory  
 167 feasibility and goal satisfaction. Since each  $\phi$  can be grounded into different action sequences, we  
 168 restrict the number of actions per subgoal to generate a finite set of possible action sequences  $\bar{a}_i$   
 169 satisfying  $\phi$ . Then, we compute the metrics for each  $\bar{a}_i$  and report the maximum score across all  $\bar{a}_i$ 's  
 170 as the trajectory feasibility and the goal satisfaction scores for  $\phi$ .

#### 171 2.6 Ability Module 4: Transition Modeling $\mathcal{T} : \langle s_0, g \rangle, o \rightarrow \langle \text{pre}, \text{eff} \rangle$

172 **Input-Output Specification.** The *transition modeling* module takes the task  $\langle s_0, g \rangle$  and a set of  
 173 operator definitions  $\{o_i\}$  as input, and generates a PDDL operator definition [8] for each  $o_i$ . In this  
 174 module, we aim to create a formal definition of actions in order to generate plans to solve the task.  
 175 During evaluation, we first extract relevant operator definitions,  $\{o_i\}$ , based on the ground truth  
 176 action trajectory  $\bar{a}$  associated with each task, with details provided in Appendix ???. Then, the LLM  
 177 generates the preconditions and effects  $\{\langle \text{pre}_i, \text{eff}_i \rangle\}$  for all operators  $\{o_i\}$ .

178 **Evaluation Metric.** The *transition modeling* module can be evaluated in two ways. First, the *logic  
 179 matching score* for an operator  $o_i$  compares the generated  $\text{pre}_i$  and  $\text{eff}_i$  against the ground truth  
 180 operator definition annotated by human experts. This comparison uses a surface form matching score  
 181 to produce an  $F_1$ -based score between two logic formulas. Intuitively, when both the LLM-generated  
 182  $\text{pre}_i$  and ground truth  $\text{pre}_i^{gt}$  are conjunctions of propositions, the  $F_1$  score is computed as the set  
 183 matching score between the sets of propositions. More complex logic formulas (e.g.,  $\forall x. \phi(x)$ ) are  
 184 evaluated recursively, as detailed in Appendix ???. The evaluation of effects is performed similarly.

185 Furthermore, the *planning success rate* assesses whether the preconditions and effects of different  
 186 operators enable a viable plan. This is computed by running an external PDDL planner [9] based

187 on generated operator definitions to achieve  $g$  from the initial state  $s_0$ . For simplicity, we only state  
188 goals in  $g$  (and ignore action subgoals). The planning success rate is 1 if the planner finds a plan.

### 189 3 Dataset Annotations and Benchmark Implementations

190 **Annotations.** Focusing on complex long-horizon tasks, we select BEHAVIOR (B) and Vir-  
191 tualHome (V) as our evaluation simulators based on their task length and scene complexity.  
192 We include a comparison of different simulators and detailed selection considerations in Ap-  
193 pendix M.1. Table 2 shows our annotations. Apart from the goal and trajectory annotations,  
194 we introduce the Goal Action annotation to reflect necessary actions that do not have post ef-  
195 fects, such as the goal action *touch* in the task “*pet the cat*”, as detailed in Appendix P.2. In  
196 the subset of VirtualHome tasks we work on, 80.7% task categories include instructions with  
197 action steps longer than 10, and 33% of the instructions have step lengths of more than 10.

199 We select BEHAVIOR as another simulator  
200 for our evaluation due to its task complexity.  
201 BEHAVIOR BDDL goals may contain quan-  
202 tifiers, such as `(forpairs (?jar ?apple)`  
203 `(inside ?apple ?jar))`, which need to  
204 be translated into grounded goals with only  
205 atomic propositions, e.g., and `((inside`  
206 `apple_1 jar_1) (inside apple_2`  
207 `jar_2))`. There can be different grounded  
208 goals that satisfy the same BDDL goal, such  
209 as `((inside apple_2 jar_1) (inside`  
210 `apple_1 jar_2))`. We call them goal options.  
211 In general, one BDDL goal corresponds to a  
212 number of goal options. The average number  
213 of grounded goals for each task is 6.7, and  
214 there are 4,164.4 goal options for each task on  
215 average. We show data distributions of goal  
216 options and other statistics in Appendix P.1.

217 **Implementation on simulators.** As BEHAVIOR does not have an action transition model layer,  
218 we implemented a symbolic simulator with an action transition model layer. Our implementation,  
219 EvalGibson, offers 30 actions that agents can use to change the states of objects. Implementation  
220 details are in Appendix Q.1. We also revise the VirtualHome simulator to support accurate evaluation,  
221 as detailed in Appendix Q.2. Evaluation settings for each large model are detailed in Appendix Q.3.

## 222 4 Results

223 We evaluate 15 open-source and proprietary LLMs on four embodied agent ability modules across  
224 two benchmark simulators: BEHAVIOR and VirtualHome. Table 3 gives an overview. Table 4,  
225 Table 5, Table 6, and Table 7 break down the analysis of four representative LLMs on four ability  
226 modules. Figure 3 shows examples of different types of error. We start with the overall analysis.

227 **Model Comparison.** Shown in Figure 2, the top performing models overall are Gemini 1.5 Pro and  
228 GPT-4o, with GPT-4o leading in its **spatial reasoning** ability and Gemini 1.5 Pro leading in its **object**  
229 **state reasoning** ability. Among all open-source models, the best performing models are Llama-3-70B  
230 and Mistral-Large-2402, while there is still a performance gap with commercial models.

231 **Ability Comparison.** GPT-4 excels in action sequencing, while Claude-3-Opus in subgoal decompo-  
232 sition and transition modeling. For open-source LLMs, Mixtral-8x22B is impressive in transition  
233 modeling and Llama-3-70B Instruct in goal interpretation. We also observe a performance gap  
234 between different simulators. Models achieve significantly lower trajectory feasibility scores on  
235 BEHAVIOR compared to VirtualHome, but achieve higher scores on goal interpretation. This is  
236 because BEHAVIOR tasks have a much longer horizon (avg 14.6 steps) while VirtualHome goals  
237 have a larger state space to search (such as “*work*”), as detailed in Appendix M.4. It shows the inverse

Table 2: Simulator dataset statistics. New annotations collected in this paper are highlighted in color.

	VirtualHome	BEHAVIOR
#task name	26	100
#task instruction	338	100
#goal	801	673
- #state	340	153
- #relation	299	520
- #action	162	-
#trajectory	338	100
- #step	2960	1460
- avg. step	8.76	14.6
#transition model	33	30
- #precondition	99	84
- #effect	57	51

Table 3: Results (%) overview. *V*: VirtualHome, *B*: BEHAVIOR. Full results in Appendix H.

Model	Goal Interpretation		Action Sequencing			Subgoal Decomposition			Transition Modeling					
	$F_1$		<i>Goal SR</i>	<i>Execution SR</i>	<i>Goal SR</i>	<i>Execution SR</i>	$F_1$	<i>Planner SR</i>	<i>V</i>	<i>B</i>	<i>V</i>	<i>B</i>		
	<i>V</i>	<i>B</i>	<i>V</i>	<i>B</i>	<i>V</i>	<i>B</i>	<i>V</i>	<i>B</i>	<i>V</i>	<i>B</i>	<i>V</i>	<i>B</i>		
Claude-3 Haiku	28.0	52.5	43.6	26.0	51.5	32.0	78.4	29.0	82.8	35.0	42.3	89.5	30.4	82.0
Claude-3 Sonnet	29.4	69.4	65.2	44.0	68.9	57.0	83.1	38.0	86.4	43.0	41.2	92.1	13.2	78.0
Claude-3 Opus	31.4	77.0	65.9	<b>51.0</b>	64.9	<b>59.0</b>	87.0	39.0	89.9	47.0	<b>48.8</b>	92.9	61.8	<b>97.0</b>
Cohere Command R	36.7	36.0	15.7	16.0	19.7	19.0	71.3	15.0	79.6	25.0	11.7	57.9	51.1	33.0
Cohere Command R+	22.4	51.2	54.8	27.0	61.0	35.0	79.0	24.0	83.7	37.0	30.8	72.1	37.2	34.0
Gemini 1.0 Pro	23.8	60.0	33.1	27.0	36.7	32.0	70.4	23.0	84.6	33.0	41.8	85.0	11.8	38.0
Gemini 1.5 Flash	26.8	74.8	61.0	40.0	65.9	52.0	<b>89.1</b>	34.0	<b>94.1</b>	42.0	45.7	86.7	46.6	75.0
Gemini 1.5 Pro	<b>37.9</b>	<b>79.6</b>	75.1	42.0	82.0	54.0	87.0	31.0	91.1	37.0	34.1	88.6	<b>91.9</b>	67.0
GPT-3.5-turbo	22.7	50.4	18.0	16.0	21.3	20.0	69.2	24.0	81.4	36.0	30.0	76.1	0.7	39.0
GPT-4-turbo	33.2	77.2	60.7	38.0	64.6	45.0	85.5	37.0	<b>94.1</b>	47.0	42.9	77.2	56.1	84.0
GPT-4o	36.5	79.2	70.2	47.0	71.8	53.0	88.8	<b>48.0</b>	90.2	<b>55.0</b>	46.7	<b>93.4</b>	68.2	90.0
Llama3 8B Instruct	22.6	28.3	22.3	10.0	22.3	16.0	48.8	21.0	58.0	29.0	12.9	66.4	28.7	47.0
Llama3 70B Instruct	26.9	70.9	54.4	34.0	56.1	42.0	78.4	20.0	87.3	30.0	37.4	80.0	12.2	12.0
Mistral Large	26.8	74.3	<b>78.4</b>	33.0	<b>83.9</b>	50.0	84.3	30.0	92.0	38.0	36.1	80.3	31.1	32.0
Mixtral 8x22B MoE	26.6	54.7	46.2	30.0	50.2	40.0	80.5	27.0	90.2	33.0	42.0	88.2	37.5	52.0

Table 4: Logic form accuracy for *goal interpretation* (%). Full results in Appendix ??.

Model	State Goal			Relation Goal			Action Goal			Overall				
	Precision	Recall	$F_1$	Precision	Recall	$F_1$	Precision	Recall	$F_1$	Precision	Recall	$F_1$		
	<i>V</i>	<i>B</i>	<i>V</i>	<i>B</i>	<i>V</i>	<i>B</i>	<i>V</i>	<i>B</i>	<i>V</i>	<i>B</i>	<i>V</i>	<i>B</i>		
Claude-3 Opus	27.0	72.6	<b>66.9</b>	93.5	38.5	81.7	22.6	75.2	<b>46.8</b>	<b>79.2</b>	30.5	77.2	14.5	-
Gemini 1.5 Pro	<b>45.4</b>	<b>94.0</b>	49.1	92.8	<b>47.2</b>	<b>93.4</b>	<b>40.0</b>	74.4	9.7	76.7	15.6	75.6	<b>26.8</b>	-
GPT-4o	29.0	67.1	60.0	94.8	39.1	78.6	31.5	<b>81.1</b>	43.6	78.5	<b>36.6</b>	<b>79.8</b>	20.5	-
Llama3 70B	23.9	69.5	61.2	<b>95.4</b>	34.3	80.4	22.6	70.0	37.5	73.3	28.2	71.6	11.2	-

correlation between trajectory evaluation performance and sequence length, as well as between goal evaluation performance and environment complexity. We further perform a systematic analysis to discover the cofactors for the goal success rate, including the number of task goals, particularly node goals, the ground truth action length, and the task object length, with details in Appendix ??.

**Object States vs Relationship.** Relational goals are generally harder to reason about compared to object-state goals. Spatial relations have a significantly lower recall in the goal interpretation task (Table 4) and a lower goal satisfaction rate (Table 5). Some non-spatial relations (e.g., *hold*) are even more difficult for LLM to predict than spatial relations, as shown in the transition modeling accuracy (Table 6): for example *holding*(toothbrush) should be a precondition for brushing teeth.

**Reporting Bias and Imprecise Physical Expressions.** Given the task “serve a meal”, all LLMs predict the incorrect goal *ontop(chicken, table)* instead of *ontop(chicken, plate)*, due to the commonly used natural language expression “put the chicken on the table”. Also, for the task “cleaning sneakers”, the goal state *onfloor(gym\_shoe, floor)* is missing from all LLM predictions, as chat models ignore the *onfloor* spatial relationship as implicit for conversational language. However, such precise physical relationships are essential for embodied task planning.

#### 4.1 Ability Module Analysis

**Goal Interpretation.** LLMs generally have difficulties distinguishing intermediate subgoals and final goals. For example, in the VirtualHome task *Drink*, GPT-4o predicts some intermediate states as part of the final goal (e.g., *open(freezer)* and *inside(water, glass)*). Overall, we observe that LLMs tend to translate NL goals word-by-word into their symbolic correspondence, rather than grounding them in the environment state. More analyses are in Appendix ??.

**Subgoal Decomposition and Action Sequencing on Trajectory Feasibility.** Most errors are runtime errors (rather than syntax errors). We illustrate examples in Figure 3. Overall, LLMs are more likely to make missing-step and additional-step errors than wrong-order or affordance errors. Missing-step errors occur when a precondition is not satisfied before the execution of an action (e.g., fetching an object without opening the box containing it). Additional steps form the most frequent errors, even for the most powerful models—it occurs when a goal has already been achieved but the

Grammar Error		Goal Satisfaction Error	
Parsing	PLACE_ONFLOOR(floor.0) ✖ Unknown action PLACE_ONFLOOR	Missing State	Missing Relation
Action-Arg Len	GRASP(rag.0, bowl.1) ✖ GRASP only has one param	LLM Output ... FIND(television.410) SWITCH_ON(television.410)	LLM Output ... next_to(plywood.78, plywood.79) and next_to(plywood.79, plywood.80)
Hallucination	RINSE(hand.65) ✖ hand.65 is not in the scene	Error Info: State Unsatisfied ✖ Missing Final State facing(agent.65, television)	Error Info: Relation Unsatisfied ✖ Missing Final Relation next_to(plywood.78, plywood.79)
Trajectory – Runtime Error			
Wrong Order	WALK(table.355) SIT(chair.356) FIND(novel.1000) GRAB(novel.1000) 	Missing Step ... CLOSE(fridge.0) SLICE(strawberry.0) SLICE(peach.0) 	Affordance Error LEFT_RELEASE OPEN(shelf.16) LEFT_RELEASE LEFT_GRASP(pool.50) 
	✖ Precondition not sitting(agent.65) = False ✓ Historical State not sitting(agent.65) = True	✖ Precondition holding(knife.0) = False ✖ Historical State holding(knife.0) = False	✖ Precondition shelf.16 not openable ✖ Precondition pool.50 not grabbable
			✖ Current State open(top_cabinet.27) = True ✖ Expected State open(top_cabinet.27) = False

Figure 3: Examples of different types of errors in trajectory feasibility, logic form parsing (e.g., in subgoals decomposition and transition modeling), and goal satisfaction rates.

Table 5: Goal satisfaction rates (%) for *action sequencing* and *subgoal decomposition*. Full results in Appendix H.2. Behavior does not include action goals.

Model	Action Sequencing						Subgoal Decomposition											
	State Goal		Relation Goal		Action Goal		Total		State Goal		Relation Goal		Action Goal		Total			
	V	B	V	B	V	B	V	B	V	B	V	B	V	B	V	B	V	B
Claude-3 Opus	63.7	45.0	71.1	53.0	77.0	-	69.1	50.8	92.4	43.0	88.6	41.6	83.3	-	89.1	42.0		
Gemini 1.5 Pro	81.7	41.0	76.7	43.2	89.2	-	82.0	42.6	91.2	31.0	72.5	37.1	89.5	-	83.9	35.4		
GPT-4o	83.1	49.0	71.1	45.5	89.9	-	81.2	46.5	92.1	50.0	84.2	53.2	93.2	-	89.4	52.3		
Llama3 70B	42.8	31.0	61.1	45.5	75.0	-	56.1	41.5	93.2	25.0	63.4	27.7	82.7	-	80.0	27.0		

Table 6: Trajectory evaluation results (%) for *action sequencing* and *subgoal decomposition*. Full results in Appendix H.3.

Model	Goal Evaluation		Trajectory Evaluation															
	Goal SR		Execution SR		Grammar Error (↓)				Runtime Error (↓)									
	V	B	V	B	V	B	V	B	V	B	V	B	V	B	V	B		
Action Sequencing																		
Claude-3 Opus	65.9	51.0	64.9	59.0	0.0	0.0	14.1	0.0	0.0	0.0	1.3	3.0	19.0	35.0	0.7	3.0	1.3	2.0
Gemini 1.5 Pro	75.1	42.0	82.0	54.0	0.3	0.0	1.6	0.0	0.3	0.0	0.0	6.0	14.8	39.0	1.0	1.0	0.7	2.0
GPT-4o	65.2	47.0	71.8	53.0	0.0	0.0	1.3	1.0	0.7	0.0	0.0	9.0	25.3	36.0	1.0	1.0	0.3	0.0
Llama3 70B	54.4	34.0	56.1	42.0	0.0	0.0	23.3	2.0	1.0	0.0	0.7	15.0	16.4	38.0	2.6	3.0	3.6	6.0
Subgoal Decomposition																		
Claude-3 Opus	87.0	39.0	89.9	47.0	0.3	0.0	3.3	3.0	0.0	0.0	1.2	5.0	3.0	45.0	2.4	0.0	16.0	5.0
Gemini 1.5 Pro	87.0	31.0	91.1	37.0	0.0	1.0	1.5	0.0	1.8	1.0	0.0	3.0	5.6	59.0	0.0	0.0	16.0	2.0
GPT-4o	88.8	48.0	90.2	55.0	0.0	0.0	6.2	3.0	0.0	0.0	1.2	5.0	2.4	37.0	0.0	0.0	15.7	5.0
Llama3 70B	78.4	20.0	87.3	30.0	0.0	1.0	2.4	5.0	0.9	1.0	2.4	8.0	5.3	51.0	1.8	4.0	20.4	4.0

model still predict to execute an additional action to achieve it (e.g., opening a box twice). More analysis is in Appendix H.3.

**Subgoal Decomposition and Action Sequencing on Goal Satisfaction Rates.** Shown in Table 5, object goals (such as *toggled\_on*) are generally easier to achieve than relational goals (such as *ontop*(agent, chair)). More analysis is provided in Appendix H.2.

**Transition Modeling.** Table 7 shows the overall performance of the logic form accuracy. For a systematic evaluation, we further categorize the tasks into five distinct ability categories requiring the transition modeling for different types of object states and relations (see Appendix I.3). Overall, we reveal significant variations in performance across different models; relational preconditions and effects are generally harder to predict than object-state ones. For instance, the Claude-3 Opus model

Table 7: Logic form accuracy ( $F_1$ ) and planner success rate (SR) for *transition modeling* (%). Full results in Appendix H.4.

Model	Object States		Object Orientation		Object Affordance		Spatial Relations		Non-Spatial Relations											
	$F_1$	SR	$F_1$	SR	$F_1$	SR	$F_1$	SR	$F_1$	SR										
	V	B	V	B	V	B	V	B	V	B										
Claude-3 Opus	<b>63.0</b>	<b>95.3</b>	63.5	<b>95.6</b>	62.6	-	71.4	-	75.5	-	58.7	-	38.7	93.1	64.8	<b>96.8</b>	7.0	<b>91.6</b>	55.4	<b>98.4</b>
Gemini 1.5 Pro	18.8	90.8	<b>94.4</b>	77.8	<b>90.9</b>	-	<b>89.3</b>	-	<b>77.7</b>	-	<b>95.8</b>	-	38.7	88.9	<b>89.0</b>	64.9	<b>7.8</b>	87.1	<b>83.8</b>	62.3
GPT-4o	54.6	93.2	71.9	91.1	52.8	-	78.6	-	74.9	-	63.5	-	<b>40.8</b>	<b>94.8</b>	66.9	90.4	7.5	<b>91.6</b>	68.9	88.5
Llama-3 70B	40.9	92.9	10.1	17.8	34.4	-	3.6	-	69.3	-	6.6	-	33.3	76.2	15.2	10.6	5.8	76.2	18.9	9.8

275 excelled in object states (63% on VirtualHome), but its performance in spatial relations is weak.  
276 Additionally, in tasks that focus on object properties, models generally perform poorly in reasoning  
277 about object orientation (e.g., the agent should be facing the TV to watch it). We also provide a  
278 sensitivity analysis tool to visualize how different transition modeling errors result in downstream  
279 planning errors (see Appendix I and H.4). We found that LLMs tend to overstate object states in  
280 effects while understating them in preconditions. Conversely, they overstate spatial relationships in  
281 preconditions and underestimate them in effects. As a result, in many cases, even if the downstream  
282 planner successfully generates a plan, it may not be feasible in the actual environment.

283 **Implications in Embodied Agent System Design.** We further investigate the potential integration of  
284 LLM-based ability modules and their robustness through **sensitive analysis** (Appendix I), **modular-**  
285 **ized vs pipeline-based** experiments (Appendix J), and **replanning** (Appendix K). We observe that  
286 trajectory feasibilities are similar, although with error accumulation from different module composi-  
287 tions, showing the potential of module composition. We have also compared different **prompting**  
288 **strategies** for embodied decision-making tasks, and summarize the best practices in Appendix L.

## 289 5 Related Work

290 Recent work in embodied decision making has been using LLMs to perform various tasks, and we  
291 include a comprehensive summary in Appendix O, see also Table 8 for a quick summary. LLMs can  
292 also be used to combine multiple of the above modules at once via chain-of-thought prompting or  
293 pipelined queries, such as goal interpretation with action sequencing [10–28], goal interpretation with  
294 subgoal decomposition [4, 23, 29], action sequencing with subgoal decomposition [23, 30, 15, 31],  
295 action sequencing with transition modeling [6, 24, 28, 32, 33, 10, 34]. Our work aims to standardize  
296 the interface between LLMs and various decision-making modules to support the seamless integration,  
297 modular evaluation, and fine-grained metrics, aiming to provide implications on using LLMs in  
298 embodied decision making more effectively and selectively. We provide additional related work on  
299 agent interfaces [35–39, 15, 40, 38, 41] and simulation benchmarks in Appendix O.

Table 8: Existing work in leveraging LLMs for embodied agents.

Goal Interpretation	Subgoal Decomposition	Action Sequencing	Transition Modeling
[4, 5, 42, 43, 18, 44, 19–28, 10–12, 45–49]	[4, 23, 30, 15, 29, 31, 50?]	[3, 6, 31, 51–54, 13, 39, 55, 14, 16, 56, 17, 38, 57, 11, 58–62, 12, 41, 63–65]	[6, 24, 28, 32, 66, 33, 10, 34]

## 300 6 Conclusions and Future Work

301 We propose a systematic evaluation framework EMBODIED AGENT INTERFACE to benchmark LLMs  
302 for embodied decision making. It focuses on 1) standardizing goal specifications using LTL formulas,  
303 2) unifying decision-making tasks through a standard interface and four fundamental ability modules,  
304 and 3) providing comprehensive fine-grained evaluation metrics and automatic error identification.  
305 We highlight the limitations of current LLMs in interpreting complex goals and different errors in  
306 reasoning, further attributing errors to various cofactors, including trajectory length, goal complexity,  
307 spatial relation goals, etc.

308 **Limitations and future work:** Our current evaluation is limited to states, actions, and goals that  
309 can be described in abstract language terms, with the input environment abstracted by relational  
310 graphs of objects. Future work should extend this to include sensory inputs and actuation outputs,  
311 possibly by extending the studied model class to include Vision-Language Models (VLMs). Other  
312 aspects of extension include the integration of memory systems (episodic memory and state memory),  
313 geometric reasoning, and navigation.

314 **References**

- 315 [1] Sanjana Srivastava, Chengshu Li, Michael Lingelbach, Roberto Martín-Martín, Fei Xia, Kent Elliott Vainio,  
316 Zheng Lian, Cem Gokmen, Shyamal Buch, Karen Liu, et al. Behavior: Benchmark for everyday household  
317 activities in virtual, interactive, and ecological environments. In *Conference on robot learning*, pages  
318 477–490. PMLR, 2022.
- 319 [2] Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Torralba.  
320 Virtualhome: Simulating household activities via programs. In *Proceedings of the IEEE conference on*  
321 *computer vision and pattern recognition*, pages 8494–8502, 2018.
- 322 [3] Jacky Liang, Wenlong Huang, F. Xia, Peng Xu, Karol Hausman, Brian Ichter, Peter R. Florence, and  
323 Andy Zeng. Code as policies: Language model programs for embodied control. *2023 IEEE International*  
324 *Conference on Robotics and Automation (ICRA)*, pages 9493–9500, 2022.
- 325 [4] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn,  
326 Keerthana Gopalakrishnan, Karol Hausman, Alexander Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz,  
327 Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, Nikhil Jayant  
328 Joshi, Ryan C. Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, Sergey Levine, Yao Lu,  
329 Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, Diego M  
330 Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, F. Xia, Ted  
331 Xiao, Peng Xu, Sichun Xu, and Mengyuan Yan. Do as i can, not as i say: Grounding language in robotic  
332 affordances. In *Conference on Robot Learning*, 2022.
- 333 [5] B. Liu, Yuqian Jiang, Xiaohan Zhang, Qian Liu, Shiqi Zhang, Joydeep Biswas, and Peter Stone. Llm+p:  
334 Empowering large language models with optimal planning proficiency. *ArXiv*, abs/2304.11477, 2023.
- 335 [6] Li Siang Wong, Jiayuan Mao, Pratyusha Sharma, Zachary S. Siegel, Jiahai Feng, Noa Korneev, Joshua B  
336 Tenenbaum, and Jacob Andreas. Learning adaptive planning representations with natural language guidance.  
337 *ArXiv*, abs/2312.08566, 2023.
- 338 [7] Jae-Woo Choi, Youngwoo Yoon, Hyobin Ong, Jaehong Kim, and Minsu Jang. Lota-bench: Benchmarking  
339 language-oriented task planners for embodied agents. *arXiv preprint arXiv:2402.08178*, 2024.
- 340 [8] Richard E Fikes and Nils J Nilsson. STRIPS: A New Approach to the Application of Theorem Proving to  
341 Problem Solving. *Artificial Intelligence*, 2(3-4):189–208, 1971.
- 342 [9] Malte Helmert. The fast downward planning system. *Journal of Artificial Intelligence Research*, 26:191–  
343 246, 2006.
- 344 [10] Huaxiaoyue Wang, Gonzalo Gonzalez-Pumariega, Yash Sharma, and Sanjiban Choudhury. Demo2code:  
345 From summarizing demonstrations to synthesizing code via extended chain-of-thought. *ArXiv*,  
346 abs/2305.16744, 2023.
- 347 [11] Boyi Li, Philipp Wu, Pieter Abbeel, and Jitendra Malik. Interactive task planning with language models.  
348 *ArXiv*, abs/2310.10645, 2023.
- 349 [12] Krishan Rana, Jesse Haviland, Sourav Garg, Jad Abou-Chakra, Ian D. Reid, and Niko Sünderhauf. Sayplan:  
350 Grounding large language models using 3d scene graphs for scalable task planning. In *Conference on*  
351 *Robot Learning*, 2023.
- 352 [13] Mengdi Xu, Peide Huang, Wenhao Yu, Shiqi Liu, Xilun Zhang, Yaru Niu, Tingnan Zhang, Fei Xia, Jie Tan,  
353 and Ding Zhao. Creative robot tool use with large language models. *ArXiv*, abs/2310.13065, 2023.
- 354 [14] Yuchen Liu, Luigi Palmieri, Sebastian Koch, Ilche Georgievski, and Marco Aiello. Delta: Decomposed  
355 efficient long-term robot task planning using large language models. *ArXiv*, abs/2404.03275, 2024.
- 356 [15] Yongchao Chen, Jacob Arkin, Yang Zhang, Nicholas A. Roy, and Chuchu Fan. Autotamp: Autoregressive  
357 task and motion planning with llms as translators and checkers. *ArXiv*, abs/2306.06531, 2023.
- 358 [16] Zhe Ni, Xiao-Xin Deng, Cong Tai, Xin-Yue Zhu, Xiang Wu, Y. Liu, and Long Zeng. Grid: Scene-graph-  
359 based instruction-driven robotic task planning. *ArXiv*, abs/2309.07726, 2023.
- 360 [17] Yike Wu, Jiatao Zhang, Nan Hu, LanLing Tang, Guilin Qi, Jun Shao, Jie Ren, and Wei Song. Mldt:  
361 Multi-level decomposition for complex long-horizon robotic task planning with open-source large language  
362 model. *ArXiv*, abs/2403.18760, 2024.

- 363 [18] Wenlong Huang, F. Xia, Ted Xiao, Harris Chan, Jacky Liang, Peter R. Florence, Andy Zeng, Jonathan  
 364 Tompson, Igor Mordatch, Yevgen Chebotar, Pierre Sermanet, Noah Brown, Tomas Jackson, Linda Luu,  
 365 Sergey Levine, Karol Hausman, and Brian Ichter. Inner monologue: Embodied reasoning through planning  
 366 with language models. In *Conference on Robot Learning*, 2022.
- 367 [19] Rishi Hazra, Pedro Zuidberg Dos Martires, and Luc De Raedt. Saycanpay: Heuristic planning with large  
 368 language models using learnable domain knowledge. *ArXiv*, abs/2308.12682, 2023.
- 369 [20] Frank Joublin, Antonello Ceravola, Pavel Smirnov, Felix Ocker, Joerg Deigmoeller, Anna Belardinelli,  
 370 Chao Wang, Stephan Hasler, Daniel Tanneberg, and Michael Gienger. Copal: Corrective planning of robot  
 371 actions with large language models. *ArXiv*, abs/2310.07263, 2023.
- 372 [21] Lihan Zha, Yuchen Cui, Li-Heng Lin, Minae Kwon, Montse Gonzalez Arenas, Andy Zeng, Fei Xia, and  
 373 Dorsa Sadigh. Distilling and retrieving generalizable knowledge for robot manipulation via language  
 374 corrections. *ArXiv*, abs/2311.10678, 2023.
- 375 [22] Wenlong Huang, Fei Xia, Dhruv Shah, Danny Driess, Andy Zeng, Yao Lu, Pete Florence, Igor Mordatch,  
 376 Sergey Levine, Karol Hausman, and Brian Ichter. Grounded decoding: Guiding text generation with  
 377 grounded models for embodied agents. In *Neural Information Processing Systems*, 2023.
- 378 [23] Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei  
 379 Lu, Xiaogang Wang, Y. Qiao, Zhaoxiang Zhang, and Jifeng Dai. Ghost in the minecraft: Generally capable  
 380 agents for open-world environments via large language models with text-based knowledge and memory.  
 381 *ArXiv*, abs/2305.17144, 2023.
- 382 [24] L. Guan, Karthik Valmeekam, Sarath Sreedharan, and Subbarao Kambhampati. Leveraging pre-trained  
 383 large language models to construct and utilize world models for model-based task planning. *ArXiv*,  
 384 abs/2305.14909, 2023.
- 385 [25] Naoki Wake, Atsushi Kanehira, Kazuhiro Sasabuchi, Jun Takamatsu, and Katsushi Ikeuchi. Chatgpt  
 386 empowered long-step robot control in various environments: A case application. *IEEE Access*, 11:95060–  
 387 95078, 2023.
- 388 [26] Zhenyu Wu, Ziwei Wang, Xiuwei Xu, Jiwen Lu, and Haibin Yan. Embodied task planning with large  
 389 language models. *ArXiv*, abs/2307.01848, 2023.
- 390 [27] Shu Wang, Muzhi Han, Ziyuan Jiao, Zeyu Zhang, Yingnian Wu, Song-Chun Zhu, and Hangxin Liu.  
 391 Llm3: Large language model-based task and motion planning with motion failure reasoning. *ArXiv*,  
 392 abs/2403.11552, 2024.
- 393 [28] Pavel Smirnov, Frank Joublin, Antonello Ceravola, and Michael Gienger. Generating consistent pddl  
 394 domains with large language models. *ArXiv*, abs/2404.07751, 2024.
- 395 [29] Chan Hee Song, Jiaman Wu, Clay Washington, Brian M. Sadler, Wei-Lun Chao, and Yu Su. Llm-  
 396 planner: Few-shot grounded planning for embodied agents with large language models. *2023 IEEE/CVF  
 397 International Conference on Computer Vision (ICCV)*, pages 2986–2997, 2022.
- 398 [30] Kolby Nottingham, Prithviraj Ammanabrolu, Alane Suhr, Yejin Choi, Hannaneh Hajishirzi, Sameer Singh,  
 399 and Roy Fox. Do embodied agents dream of pixelated sheep?: Embodied decision making using language  
 400 guided world modelling. In *International Conference on Machine Learning*, 2023.
- 401 [31] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi (Jim) Fan,  
 402 and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *ArXiv*,  
 403 abs/2305.16291, 2023.
- 404 [32] S. Sundar Raman, Vanya Cohen, Eric Rosen, Ifrah Idrees, David Paulius, and Stefanie Tellex. Cape:  
 405 Corrective actions from precondition errors using large language models. In *ICRA*, 2022.
- 406 [33] Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter  
 407 Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using large  
 408 language models. *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages  
 409 11523–11530, 2022.
- 410 [34] Zhaoyi Li, Kelin Yu, Shuo Cheng, and Danfei Xu. LEAGUE++: EMPOWERING CONTINUAL ROBOT  
 411 LEARNING THROUGH GUIDED SKILL ACQUISITION WITH LARGE LANGUAGE MODELS. In  
 412 *ICLR 2024 Workshop on Large Language Model (LLM) Agents*, 2024.

- 413 [35] Georgios Fainekos, Hadas Kress-Gazit, and George Pappas. Temporal logic motion planning for mobile  
 414 robots. *Proceedings of the 2005 IEEE International Conference on Robotics and Automation*, pages  
 415 2020–2025, 2005.
- 416 [36] Hadas Kress-Gazit, Georgios Fainekos, and George Pappas. Temporal-logic-based reactive mission and  
 417 motion planning. *IEEE Transactions on Robotics*, 25:1370–1381, 2009.
- 418 [37] Stephen L. Smith, Jana Tumova, Calin A. Belta, and Daniela Rus. Optimal path planning for surveillance  
 419 with temporal-logic constraints\*. *The International Journal of Robotics Research*, 30:1695 – 1708, 2011.
- 420 [38] A. Mavrogiannis, Christoforos Mavrogiannis, and Yiannis Aloimonos. Cook2ltl: Translating cooking  
 421 recipes to ltl formulae using large language models. *ArXiv*, abs/2310.00163, 2023.
- 422 [39] J. Wang, Jiaming Tong, Kai Liang Tan, Yevgeniy Vorobeychik, and Yiannis Kantaros. Conformal temporal  
 423 logic planning using large language models: Knowing when to do what and when to ask for help. *ArXiv*,  
 424 abs/2309.10092, 2023.
- 425 [40] Amir Pnueli. The temporal logic of programs. *18th Annual Symposium on Foundations of Computer  
 426 Science (sfcs 1977)*, pages 46–57, 1977.
- 427 [41] Wenqi Zhang, Ke Tang, Hai Wu, Mengna Wang, Yongliang Shen, Guiyang Hou, Zeqi Tan, Peng Li, Yueting  
 428 Zhuang, and Weiming Lu. Agent-pro: Learning to evolve via policy-level reflection and optimization,  
 429 2024.
- 430 [42] Yan Ding, Xiaohan Zhang, Chris Paxton, and Shiqi Zhang. Task and motion planning with large language  
 431 models for object rearrangement. *2023 IEEE/RSJ International Conference on Intelligent Robots and  
 432 Systems (IROS)*, pages 2086–2092, 2023.
- 433 [43] Yaqi Xie, Chenyao Yu, Tongyao Zhu, Jinbin Bai, Ze Gong, and Harold Soh. Translating natural language  
 434 to planning goals with large-language models. *ArXiv*, abs/2302.05128, 2023.
- 435 [44] Kevin Lin, Christopher Agia, Toki Migimatsu, Marco Pavone, and Jeannette Bohg. Text2motion: from  
 436 natural language instructions to feasible plans. *Autonomous Robots*, 47:1345 – 1365, 2023.
- 437 [45] Zeyuan Yang, Jiageng Liu, Peihao Chen, Anoop Cherian, Tim K Marks, Jonathan Le Roux, and Chuang  
 438 Gan. Rila: Reflective and imaginative language agent for zero-shot semantic audio-visual navigation.
- 439 [46] Yang Zhang, Shixin Yang, Chenjia Bai, Fei Wu, Xiu Li, Zhen Wang, and Xuelong Li. Towards efficient  
 440 llm grounding for embodied multi-agent collaboration, 2024.
- 441 [47] Xudong Guo, Kaixuan Huang, Jiale Liu, Wenhui Fan, Natalia Vélez, Qingyun Wu, Huazheng Wang,  
 442 Thomas L. Griffiths, and Mengdi Wang. Embodied llm agents learn to cooperate in organized teams, 2024.
- 443 [48] Hongxin Zhang, Weihua Du, Jiaming Shan, Qinhong Zhou, Yilun Du, Joshua B Tenenbaum, Tianmin Shu,  
 444 and Chuang Gan. Building cooperative embodied agents modularly with large language models. In *The  
 445 Twelfth International Conference on Learning Representations*, 2023.
- 446 [49] Yue Wu, Xuan Tang, Tom M. Mitchell, and Yuanzhi Li. Smartplay: A benchmark for llms as intelligent  
 447 agents, 2024.
- 448 [50] Enshen Zhou, Yiran Qin, Zhenfei Yin, Yuzhou Huang, Ruimao Zhang, Lu Sheng, Yu Qiao, and Jing Shao.  
 449 Minedreamer: Learning to follow instructions via chain-of-imagination for simulated-world control, 2024.
- 450 [51] Tom Silver, Soham Dan, Kavitha Srinivas, Joshua B. Tenenbaum, Leslie Pack Kaelbling, and Michael  
 451 Katz. Generalized planning in pddl domains with pretrained large language models. In *AAAI Conference  
 452 on Artificial Intelligence*, 2023.
- 453 [52] Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu.  
 454 Reasoning with language model is planning with world model. *ArXiv*, abs/2305.14992, 2023.
- 455 [53] Jacky Liang, Fei Xia, Wenhao Yu, Andy Zeng, Montse Gonzalez Arenas, Maria Attarian, Maria Bauza,  
 456 Matthew Bennice, Alex Bewley, Adil Dostmohamed, Chuyuan Fu, Nimrod Gileadi, Marissa Giustina,  
 457 Keerthana Gopalakrishnan, Leonard Hasenclever, Jan Humplik, Jasmine Hsu, Nikhil Joshi, Ben Jyenis,  
 458 Chase Kew, Sean Kirmani, Tsang-Wei Edward Lee, Kuang-Huei Lee, Assaf Hurwitz Michaely, Joss  
 459 Moore, Kenneth Oslund, Dushyant Rao, Allen Z. Ren, Baruch Tabanpour, Quan Ho Vuong, Ayzaan Wahid,  
 460 Ted Xiao, Ying Xu, Vincent Zhuang, Peng Xu, Erik Frey, Ken Caluwaerts, Ting-Yu Zhang, Brian Ichter,  
 461 Jonathan Tompson, Leila Takayama, Vincent Vanhoucke, Izhak Shafran, Maja Mataric, Dorsa Sadigh,  
 462 Nicolas Manfred Otto Heess, Kanishka Rao, Nik Stewart, Jie Tan, and Carolina Parada. Learning to learn  
 463 faster from human feedback with language model predictive control. *ArXiv*, abs/2402.11450, 2024.

- 464 [54] Yongchao Chen, Jacob Arkin, Yilun Hao, Yang Zhang, Nicholas Roy, and Chuchu Fan. Prompt optimization  
 465 in multi-step tasks (promst): Integrating human feedback and preference alignment. *ArXiv*, abs/2402.08702,  
 466 2024.
- 467 [55] Yingdong Hu, Fanqi Lin, Tong Zhang, Li Yi, and Yang Gao. Look before you leap: Unveiling the power  
 468 of gpt-4v in robotic vision-language planning. *ArXiv*, abs/2311.17842, 2023.
- 469 [56] Georgia Chalvatzaki, Ali Younes, Daljeet Nandha, An T. Le, Leonardo F. R. Ribeiro, and Iryna Gurevych.  
 470 Learning to reason over scene graphs: a case study of finetuning gpt-2 into a robot language model for  
 471 grounded task planning. *Frontiers in Robotics and AI*, 10, 2023.
- 472 [57] Mandi Zhao, Shreya Jain, and Shuran Song. Roco: Dialectic multi-robot collaboration with large language  
 473 models. *ArXiv*, abs/2307.04738, 2023.
- 474 [58] Huaxiaoyue Wang, K. Kedia, Juntao Ren, Rahma Abdullah, Atiksh Bhardwaj, Angela Chao, Kelly Y  
 475 Chen, Nathaniel Chin, Prithwish Dan, Xinyi Fan, Gonzalo Gonzalez-Pumariega, Aditya Kompella, Max-  
 476 imus Adrian Pace, Yash Sharma, Xiangwan Sun, Neha Sunkara, and Sanjiban Choudhury. Mosaic: A  
 477 modular system for assistive and interactive cooking. *ArXiv*, abs/2402.18796, 2024.
- 478 [59] Murtaza Dalal, Tarun Chiruvolu, Devendra Singh Chaplot, and Ruslan Salakhutdinov. Plan-seq-learn:  
 479 Language model guided rl for solving long horizon robotics tasks. In *ICLR*, 2024.
- 480 [60] Meenal Parakh, Alisha Fong, Anthony Simeonov, Abhishek Gupta, Tao Chen, and Pulkit Agrawal. Lifelong  
 481 robot learning with human assisted language planners. *arXiv:2309.14321*, 2023.
- 482 [61] Zeyi Liu, Arpit Bahety, and Shuran Song. Reflect: Summarizing robot experiences for failure explanation  
 483 and correction. *ArXiv*, abs/2306.15724, 2023.
- 484 [62] Wenlong Huang, P. Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners:  
 485 Extracting actionable knowledge for embodied agents. *ArXiv*, abs/2201.07207, 2022.
- 486 [63] Yiran Qin, Enshan Zhou, Qichang Liu, Zhenfei Yin, Lu Sheng, Ruimao Zhang, Yu Qiao, and Jing Shao.  
 487 Mp5: A multi-modal open-ended embodied system in minecraft via active perception, 2024.
- 488 [64] Qinhong Zhou, Sunli Chen, Yisong Wang, Haozhe Xu, Weihua Du, Hongxin Zhang, Yilun Du, Joshua B  
 489 Tenenbaum, and Chuang Gan. Hazard challenge: Embodied decision making in dynamically changing  
 490 environments. *arXiv preprint arXiv:2401.12975*, 2024.
- 491 [65] Zhonghan Zhao, Wenhao Chai, Xuan Wang, Li Boyi, Shengyu Hao, Shidong Cao, Tian Ye, Jenq-Neng  
 492 Hwang, and Gaoang Wang. See and think: Embodied agent in virtual environment, 2023.
- 493 [66] Yan Ding, Xiaohan Zhang, S. Amiri, Nieqing Cao, Hao Yang, Andy Kaminski, Chad Esselink, and Shiqi  
 494 Zhang. Integrating action knowledge and llms for task planning and situation handling in open worlds.  
 495 *Autonomous Robots*, 47:981 – 997, 2023.
- 496 [67] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang,  
 497 Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang,  
 498 Christian Cosgrove, Christopher D. Manning, Christopher R'e, Diana Acosta-Navas, Drew A. Hudson,  
 499 E. Zelikman, Esin Durmus, Faisal Ladhab, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav  
 500 Santhanam, Laurel J. Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun, Nathan S. Kim, Neel Guha,  
 501 Niladri S. Chatterji, O. Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani  
 502 Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas F. Icard, Tianyi Zhang, Vishrav Chaudhary,  
 503 William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. Holistic evaluation of language  
 504 models. *Annals of the New York Academy of Sciences*, 1525:140 – 146, 2023.
- 505 [68] Richard Bellman. A markovian decision process. *Indiana University Mathematics Journal*, 1957.
- 506 [69] Thomas L. Dean and Michael P. Wellman. *Planning and Control*. 1991.
- 507 [70] Tom Silver and Rohan Chitnis. Pddlgym: Gym environments from pdl problems. *ArXiv*, abs/2002.06432,  
 508 2020.
- 509 [71] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny  
 510 Zhou. Chain of thought prompting elicits reasoning in large language models. *ArXiv*, abs/2201.11903,  
 511 2022.
- 512 [72] Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select: Inter-  
 513 active planning with large language models enables open-world multi-task agents. *ArXiv*, abs/2302.01560,  
 514 2023.

- 515 [73] Marta Skreta, Zihan Zhou, Jia Lin Yuan, Kourosh Darvish, Alán Aspuru-Guzik, and Animesh Garg. Replan:  
516 Robotic replanning with perception and language models, 2024.
- 517 [74] Karthik Valmeekam, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. Planbench: An  
518 extensible benchmark for evaluating large language models on planning and reasoning about change. In  
519 *Neural Information Processing Systems*, 2022.
- 520 [75] Chengshu Li, Fei Xia, Roberto Martín-Martín, Michael Lingelbach, Sanjana Srivastava, Bokui Shen, Kent  
521 Vainio, Cem Gokmen, Gokul Dharan, Tanish Jain, Andrey Kurenkov, C. Karen Liu, Hyowon Gweon,  
522 Jiajun Wu, Li Fei-Fei, and Silvio Savarese. igibson 2.0: Object-centric simulation for robot learning of  
523 everyday household tasks, 2021.

524 **Checklist**

- 525 1. For all authors...
  - 526 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
527 contributions and scope? **[Yes]** We clearly state our problem scope and contributions.
  - 528 (b) Did you describe the limitations of your work? **[Yes]** See Appendix Section ??.
  - 529 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** See  
530 Appendix Section ??.
  - 531 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
532 them? **[Yes]**
- 533 2. If you are including theoretical results...
  - 534 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
  - 535 (b) Did you include complete proofs of all theoretical results? **[N/A]**
- 536 3. If you ran experiments (e.g. for benchmarks)...
  - 537 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
538 perimental results (either in the supplemental material or as a URL)? **[Yes]** All code,  
539 annotations, and instructions for reproducing our results are included in the supple-  
540 mentary materials.
  - 541 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
542 were chosen)? **[Yes]**
  - 543 (c) Did you report error bars (e.g., with respect to the random seed after running exper-  
544 iments multiple times)? **[N/A]** LLM inference tasks are very resource intensive and  
545 proprietary model APIs are too costly.
  - 546 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
547 of GPUs, internal cluster, or cloud provider)? **[Yes]**
- 548 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - 549 (a) If your work uses existing assets, did you cite the creators? **[Yes]** We cite the existing  
550 assets in the reference.
  - 551 (b) Did you mention the license of the assets? **[Yes]** We mention them in our released  
552 website and in the supplemental material.
  - 553 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**  
554 In the supplemental material.
  - 555 (d) Did you discuss whether and how consent was obtained from people whose data you're  
556 using/curating? **[Yes]** The resources are from existing public data that is open-access.
  - 557 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
558 information or offensive content? **[Yes]** The data we are using does not contain  
559 personally identifiable information or offensive content.
- 560 5. If you used crowdsourcing or conducted research with human subjects...
  - 561 (a) Did you include the full text of instructions given to participants and screenshots, if  
562 applicable? **[N/A]**

- 563                     (b) Did you describe any potential participant risks, with links to Institutional Review  
564                         Board (IRB) approvals, if applicable? [N/A]  
565                     (c) Did you include the estimated hourly wage paid to participants and the total amount  
566                         spent on participant compensation? [N/A]

# 568 Appendix

## 569 570 Table of Contents

---

571	<b>A Summary of Empirical Findings</b>	<b>18</b>
572	<b>B Datasheets for EMBODIED AGENT INTERFACE (EAGENT)</b>	<b>20</b>
573	<b>C URLs of Data and License</b>	<b>27</b>
574	C.1 Hosting and Maintenance Plan . . . . .	27
575	C.2 Long-term Preservation and DOI . . . . .	27
576	C.3 URL of Croissant Metadata Record . . . . .	28
577	C.4 Dataset Format . . . . .	28
578	C.5 Author Statement . . . . .	30
579	<b>D URLs of Code and Re-productivity</b>	<b>30</b>
580	D.1 Code for VirtualHome Evaluator and Computation Resources . . . . .	30
581	D.2 Code for BEHAVIOR Evaluator and Computation Resources . . . . .	30
582	D.3 Code for LLMs Implementations and Computation Resources . . . . .	30
583	<b>E Impact, Limitations and Future Directions</b>	<b>31</b>
584	E.1 Broader Impact . . . . .	31
585	E.2 Limitations . . . . .	32
586	E.3 Potential Negative Social Impact . . . . .	32
587	E.4 Potential Future Directions . . . . .	32
588	<b>F Embodied Agent Interface Details</b>	<b>33</b>
589	F.1 Input and Output Details . . . . .	33
590	F.2 Grounding to Markov Decision Process . . . . .	34
591	F.3 The Relationship between Ability Modules . . . . .	36
592	<b>G Fine-Grained Metrics and Automatic Error Detection</b>	<b>37</b>
593	G.1 Goal Interpretation: State Goal, Relation Goal and Action Goal . . . . .	37
594	G.2 Action Sequencing: Trajectory Error Detection for Missing Step, Additional Step, Wrong Temporal Order, Affordance Error . . . . .	38
595	G.3 Subgoal Decomposition: Converting Subgoal Trajectory to Action Trajectory with BFS Searching . . . . .	40
596	G.4 Transition Modeling: Evaluating with PDDL Planners . . . . .	41
599	<b>H Full Results with 15 models</b>	<b>43</b>
600	H.1 Goal Interpretation . . . . .	43
601	H.2 Subgoal Decomposition . . . . .	44
602	H.3 Action Sequencing . . . . .	49
603	H.4 Transition Modeling . . . . .	51
604	H.5 Correlation with Action Length and Goal Complexity . . . . .	57
605	<b>I Sensitivity Analysis</b>	<b>58</b>
606	I.1 Motivation and Problem Formulation . . . . .	58
607	I.2 Implementation Details . . . . .	60
608	I.3 Result Analysis . . . . .	60
609	<b>J Pipeline-Based vs Modularized</b>	<b>64</b>
610	J.1 Motivation and Problem Formulation . . . . .	64
611	J.2 Implementation Details . . . . .	64

612	J.3 Result Analysis . . . . .	64
613	<b>K Replanning and Feedback</b>	<b>65</b>
614	K.1 Motivation and Problem Formulation . . . . .	65
615	K.2 Implementation Details . . . . .	65
616	K.3 Result Analysis . . . . .	66
617	<b>L Prompt and Analysis</b>	<b>66</b>
618	L.1 Prompt of Goal Interpretation . . . . .	66
619	L.2 Prompt of Subgoal Decomposition . . . . .	67
620	L.3 Prompt of Action Sequencing . . . . .	68
621	L.4 Prompt of Transition Modeling . . . . .	72
622	L.5 Prompt of Environment Representation . . . . .	76
623	L.6 Prompt Analysis and Learned Lessons . . . . .	77
624	<b>M Dataset Analysis</b>	<b>78</b>
625	M.1 Simulator Comparison and Selection . . . . .	78
626	M.2 Data Structure . . . . .	79
627	M.3 Goal Complexity Analysis . . . . .	79
628	M.4 Data Statistics and Distribution . . . . .	80
629	M.5 Task List . . . . .	81
630	M.6 Task Categorization . . . . .	83
631	<b>N LTL Implementation Details</b>	<b>84</b>
632	N.1 Why LTL . . . . .	84
633	N.2 Syntax and Semantics of LTL Formulas . . . . .	85
634	<b>O Extensive Related Work</b>	<b>86</b>
635	O.1 LLMs for Embodied Planning . . . . .	86
636	O.2 LTL Agent Interface . . . . .	87
637	<b>P Annotation Details</b>	<b>88</b>
638	P.1 BEHAVIOR . . . . .	88
639	P.2 VirtualHome . . . . .	88
640	<b>Q Simulator Implementation Details</b>	<b>89</b>
641	Q.1 BEHAVIOR Implementation Details . . . . .	89
642	Q.2 VirtualHome Implementation Details . . . . .	91
643	Q.3 Evaluation Settings of LLMs . . . . .	92
644		
645		
646		

---

647 **A Summary of Empirical Findings**

648 **1. Goal Interpretation:**

- 649 • Most LLMs still struggle to faithfully translate natural language instructions into  
650 grounded states (objects, object states, and relations) in the environment.
- 651 • A common error is generating intermediate goals instead of final goals, e.g., predicting  
652 the state *open*(freezer) for the task “drinking water”.
- 653 • Another common error is omitting conversationally uncommon spatial relationship  
654 goals. For example, in the task “serving a meal”, with ground truth goal condition  
655 *ontop(chicken.0, plate.2)* and *ontop(plate.2, table.1)*, GPT-4o mistakenly predicts  
656 *ontop(chicken.0, table.1)*, ignoring the crucial spatial relationship between the chicken,  
657 plate, and table.
- 658 • Gemini 1.5 Pro achieves the highest overall goal interpretation performance (F1-score)  
659 in both VirtualHome and BEHAVIOR simulators, while Claude-3 Opus has the highest  
660 successful ground truth goal retrieval rate (Recall) in both simulators. For example,  
661 in the VirtualHome simulator, Gemini 1.5 Pro achieves an F1-score of 82.0%, and  
662 Claude-3 Opus achieves a Recall of 89.1%.
- 663 • State-of-the-art proprietary LLMs make few to no grammar errors, while top open-  
664 source LLMs like Llama3 70B Instruct suffer more from format/parsing errors and  
665 object/state hallucination. For instance, GPT-4o makes no parsing errors in both  
666 simulators, while Llama3 8B makes parsing errors in 0.6% of cases in VirtualHome  
667 and 2.0% in BEHAVIOR.

668 **2. Action Sequencing:**

- 669 • Reasoning ability is a crucial aspect that LLMs should improve. As show in Figure 1 in  
670 the main paper, trajectory feasibility errors are common (45.2%), with a large portion  
671 of missing step (19.5%) and additional step (14.2%) errors, often due to overlooking  
672 preconditions. For instance, LLMs may ignore the agent’s *sitting* or *lying* state and  
673 fail to include a *standup* action before executing other actions. They sometimes also  
674 fail to understand the need to *open* a *closed* object before *fetching* items from inside.  
675 Additional step errors frequently occur when LLMs output actions for previously  
676 achieved goals.
- 677 • GPT-4o achieves the highest goal success rate and execution success rate in both  
678 simulators, with a goal success rate of 81.2% and an execution success rate of 53.0%  
679 in VirtualHome, and 46.5% and 53.0% in BEHAVIOR, respectively.
- 680 • SOTA LLMs generally make fewer grammar errors compared to less advanced models.  
681 For example, Claude-3 Opus makes no parsing errors in both simulators, while GPT-  
682 3.5-turbo makes parsing errors in 4.0% of cases in BEHAVIOR.
- 683 • The most common runtime errors are missing steps and wrong order in both simulators.  
684 For instance, in BEHAVIOR, GPT-4o encounters missing step errors in 36.0% of cases  
685 and wrong order errors in 9.0% of cases.
- 686 • LLMs perform better in satisfying state goals than relation goals and struggle with  
687 complex action goals. For example, in VirtualHome, GPT-4o achieves a state goal  
688 success rate of 83.1% but a relation goal success rate of 71.1%.
- 689 • Task complexity, including the number of goals, state goals, relation goals, and action  
690 sequence length, adversely affects the goal success rate. For instance, in BEHAVIOR,  
691 the success rate drops from around 60% for tasks with fewer than 5 goals to below  
692 40% for tasks with more than 10 goals. Subgoal decomposition is not strictly easier  
693 than action sequencing in abstract action spaces. Subgoal decomposition is not strictly  
694 easier than action sequencing in abstract action spaces. Subgoal decomposition is not  
695 strictly easier than action sequencing in abstract action spaces. Subgoal decomposition  
696 is not strictly easier than action sequencing in abstract action spaces. Subgoal decom-  
697 position is not strictly easier than action sequencing in abstract action spaces. Subgoal

698 decomposition is not strictly easier than action sequencing in abstract action spaces.  
699 Subgoal decomposition is not strictly easier than action sequencing in abstract action  
700 spaces. Subgoal decomposition is not strictly easier than action sequencing in abstract  
701 action spaces.

702 **3. Subgoal Decomposition:**

- 703 • Subgoal decomposition is not strictly easier than action sequencing in abstract action  
704 spaces.
- 705 • GPT-4o demonstrates superior performance in both VirtualHome and BEHAVIOR  
706 simulators compared to other SOTA LLMs, with success rates of 88.8% and 48.0%,  
707 respectively. Gemini 1.5 Flash also shows the highest performance in VirtualHome  
708 with a success rate of 89.1%.
- 709 • SOTA models generally avoid grammar errors but can hallucinate actions and objects.  
710 For example, GPT-4o tends to hallucinate the action *POUR* when dealing with the  
711 task “make coffee” in VirtualHome, which is not defined in the subgoal decomposition  
712 setting.
- 713 • The most common runtime errors in VirtualHome are additional steps, while in BE-  
714 HAVIOR, they are missing steps. For instance, in VirtualHome, all LLMs are prone to  
715 produce additional steps errors, even for SOTA LLMs like GPT-4o and Claude-3 Opus.  
716 This is mainly because, in the initial scene state, some of the goals have already been  
717 achieved, yet LLMs still prefer to plan the satisfied goals in their output.
- 718 • Stronger LLMs like GPT-4o show higher accuracy in action goal success rates in  
719 VirtualHome compared to weaker models like Llama3 8B. However, achieving state and  
720 relation goals in BEHAVIOR is challenging due to more complex task representations  
721 and stricter precondition checks. For example, in BEHAVIOR, most state and relation  
722 goals are encapsulated within quantifiers, and quantifiers such as “forall” or “forpairs”  
723 tend to fail if even a single state or relation goal is not met.
- 724 • Overall LLM performance is lower in BEHAVIOR compared to VirtualHome due to  
725 complex task representations involving quantifiers like “forall” and “forpairs”, which  
726 articulate complex temporal and spatial requirements. For instance, most tasks in  
727 BEHAVIOR have quantifiers with complex spatial or temporal requirements, while  
728 VirtualHome tasks have much easier goal definitions.

729 **4. Transition Modeling:**

- 730 • Models like Claude-3 Opus and Gemini 1.5 Pro excel in specific categories like  
731 object states and object orientation, suggesting that targeted training or specialized  
732 architectures enhance LLM capabilities in understanding different types of tasks in  
733 transition modeling. For example, Claude-3 Opus achieves an F1-score of 95.3% in  
734 object states in BEHAVIOR, while Gemini 1.5 Pro achieves an F1-score of 90.9% in  
735 object orientation in VirtualHome.
- 736 • Across various models, non-spatial relations consistently pose a challenge, highlighting  
737 a gap in the ability of LLMs to grasp complex relational dynamics. For instance, in  
738 VirtualHome, the best-performing model, Gemini 1.5 Flash, only achieves an F1-score  
739 of 7.9% in non-spatial relations.
- 740 • The effectiveness of planning relies heavily on the consistency of the predicted action  
741 space by LLMs; discrepancies between mixed predicted and ground truth actions lead  
742 to reduced planner success. For example, if we mix the action spaces of GPT-4o pre-  
743 dictions and ground truth, using “plug\_in” from GPT-4o prediction and “walk\_toward”  
744 and “switch\_on” from ground truth, the PDDL planner cannot find a feasible solution  
745 for the task.

746 **5. Sensitivity Analysis:**

- 747 • Specific actions like “plug\_in” and “walk\_towards” consistently show low success rates  
748 due to complex preconditions and spatial requirements. For instance, in VirtualHome,  
749 the success rate for “plug\_in” is only 0.09, and for “walk\_towards”, it is 0.63.

- 750           • Complex interactions involving detailed object manipulation, such as  
 751           “slice\_carvingknife” and “place\_inside”, present notable challenges. For ex-  
 752           ample, in BEHAVIOR, the success rate for “slice\_carvingknife” is 0.00, and for  
 753           “place\_inside”, it shows a rather low success rate in many tasks.  
 754           • Current training regimens may not fully capture the diversity of real-world interactions,  
 755           especially in spatial and object-oriented tasks. This is evident from the generally lower  
 756           success rates for actions involving complex spatial relationships and object interactions.

757           **6. Pipeline-Based vs. Modularized:**

- 758           • Both modularized and pipeline-based methods have similar trajectory executable rates.  
 759           For example, in the pipeline of Goal Interpretation and Action Sequencing in BEHAV-  
 760           IOR, the modularized method has an execution success rate of 53.0% for GPT-4o,  
 761           while the pipeline-based method has an execution success rate of 55.0%.  
 762           • Pipeline-based methods suffer from error accumulation due to the composition of two  
 763           modules. For instance, in the pipeline of Goal Interpretation and Subgoal Decom-  
 764           position in BEHAVIOR, the goal success rate for GPT-4o drops from 48.0% in the  
 765           modularized method to 38.0% in the pipeline-based method.  
 766           • SOTA LLMs generally avoid grammar errors for both pipeline-based and modularized  
 767           methods, unlike less advanced models. For example, GPT-4o makes no parsing errors  
 768           in both methods, while Llama3 8B makes parsing errors in 2.0% of cases in the  
 769           pipeline-based method.  
 770           • All LLMs, regardless of their advancement, are prone to runtime errors, missing  
 771           necessary steps in their generation process. For instance, in the pipeline of Goal  
 772           Interpretation and Action Sequencing in BEHAVIOR, GPT-4o encounters missing step  
 773           errors in 35.0% of cases in both modularized and pipeline-based methods.

774           **7. Replanning and Feedback:**

- 775           • Incorporating replanning based on feedback significantly improves the model’s per-  
 776           formance, demonstrating over a 10% increase in success rates. For example, with  
 777           replanning, GPT-4o’s goal success rate increases from 47.0% to 59.0%, and its  
 778           execution success rate increases from 53.0% to 63.0% in BEHAVIOR .  
 779           • Replanning can sometimes result in the over-generation of actions, as indicated by  
 780           an increased rate of additional steps errors. For instance, with replanning, GPT-4o’s  
 781           additional step error rate increases from 0.0% to 3.0% in BEHAVIOR .

782       These empirical findings, along with the provided examples, highlight the strengths and weaknesses  
 783       of LLMs in embodied decision-making tasks across different ability modules and simulators. The  
 784       insights gained from these experiments can guide future research and development efforts to address  
 785       the identified challenges and improve the performance of LLM-based embodied agents. We present  
 786       more examples in Appendix H to illustrate the specific areas where LLMs excel or struggle, providing  
 787       a more concrete understanding of their capabilities and limitations in various scenarios.

788           **B Datasheets for EMBODIED AGENT INTERFACE (EAGENT)**

789           **Motivation**

791       **For what purpose was the dataset created?** Was there a specific task in mind? Was there a specific  
 792       gap that needed to be filled? Please provide a description.

793       EMBODIED AGENT INTERFACE (EAGENT) is a diagnostic benchmark for assessing **embodied**  
 794       **decision making** capabilities of LLMs. Given the natural language instructions from humans (such  
 795       as “*cleaning the refrigerator*”, “*polishing furniture*”), LLMs serve as embodied agents to achieve the  
 796       specified goals through a sequence of actions in various embodied environments.

797       The key difference between language models and embodied agent models is the ability of (1)  
 798       interacting with the environment, (2) goal-driven, and (3) decision making to achieve the goal, As

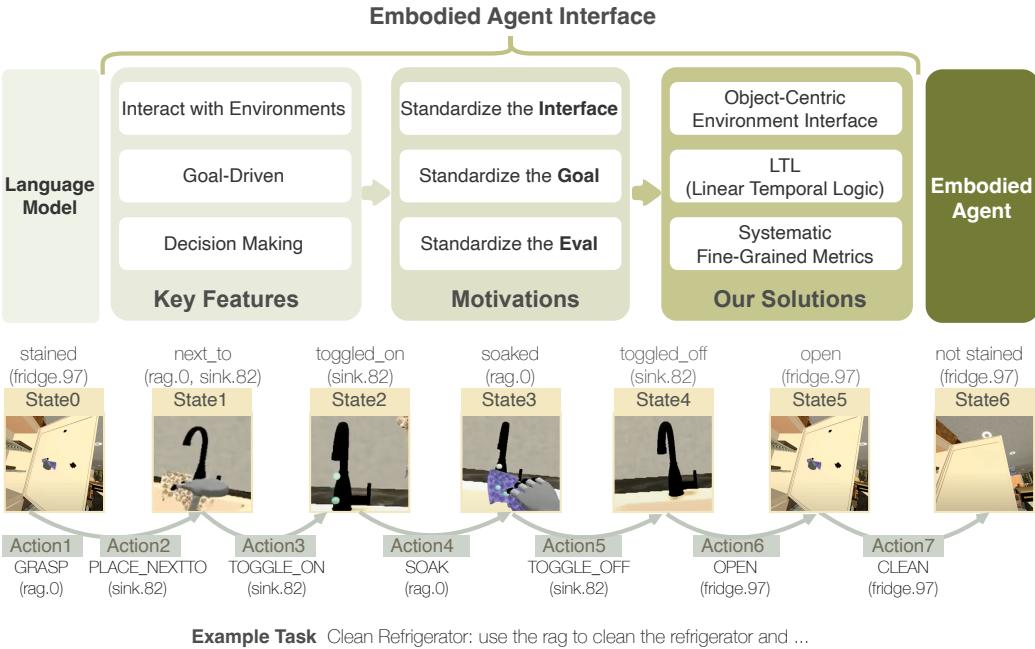


Figure 4: Compared to general language models, the embodied agent has three key new abilities, including interacting with environments, goal-driven and perform decision making to achieve the goal. We believe a systematic evaluation for embodied decision making should cover three aspects by standardizing the interface, the goal representation and the evaluation metrics. Our EMBODIED AGENT INTERFACE address the limitations of traditional evaluations by focusing on goal-driven evaluation, standard interface and modules, as well as broad coverage of evaluation and fine-grained metrics.

799 shown in Figure 4. While some prior works [7] have proposed benchmarks with simulators to  
800 validate the output plan with a success rate, they are in the high-level natural language planning  
801 space without connecting to objects and state changes in the embodied environment, as shown in  
802 Figure 5. We address the limitations of traditional evaluations on benchmarking embodied decision  
803 making from three aspects: (1) “benchmarking”: We propose a **broad coverage of evaluation and**  
804 **fine-grained metrics**. Our interface offers fine-grained metrics to automatically identify various  
805 error types (such as missing step, additional step, wrong temporal order, affordance error, etc),  
806 providing a comprehensive evaluation of LLM performance. (2) “embodied”: We move from high-  
807 level natural language planning to lower-level object interactions in the embodied environment.  
808 We **standardize goal specifications as linear temporal logic (LTL) formulas based on object-**  
809 **centric representations**, extending goals from states to temporally dependent logical transitions. (3)  
810 “decision making”: We **standardize interface and modules approach** by unifying a broad set of  
811 decision-making tasks involving states and temporally extended goals, four key LLM-based modules  
812 (goal interpretation, subgoal decomposition, action sequencing, and transition modeling), covering  
813 the fundamental abilities in the Markov Decision Process (detailed in Appendix F.2). Please see  
814 Section 1 in the main paper for more details.

815 **Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g.,**  
816 **company, institution, organization)?**

817 The authors created the dataset within the Stanford Vision and Learning Lab at Stanford University.  
818 It is created for public use and is not affiliated with or tied to any specific organization or institution.

819 **Who funded the creation of the dataset?** If there is an associated grant, please provide the name of  
820 the grantor and the grant name and number.

821 The benchmark creation is funded by the 1287101-1-UBKNA.

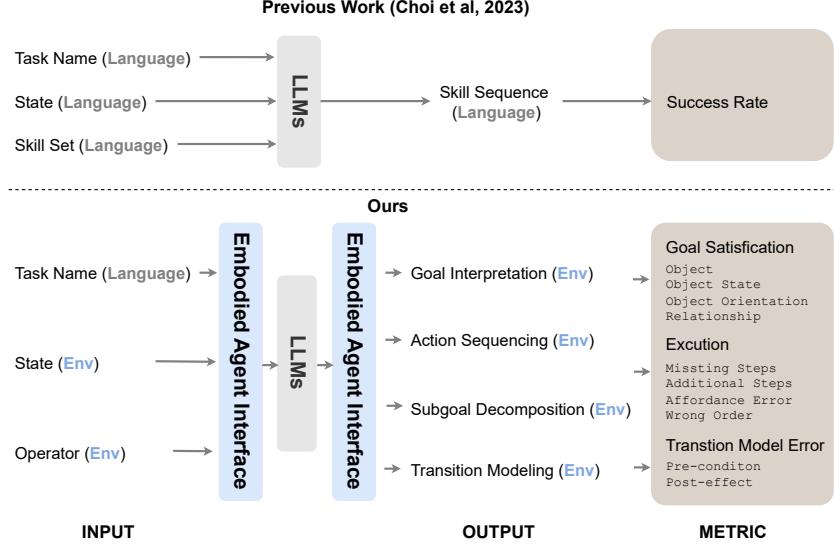


Figure 5: Comparison with existing benchmarks on LLMs for embodied decision making.

822 Any other comments?

823 No.

824

## Composition

825

826 **What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?** Are there multiple types of instances (e.g., movies, users, and ratings; people and interactions between them; nodes and edges)? Please provide a description.

829 Each instance in the dataset represents a task with its corresponding ground truth plan. Specifically, 830 each task contains the following data: (1) natural language task name, (2) natural language task 831 instruction, (3) symbolic goal definition, (4) symbolic action trajectory, (5) the transition models 832 involved in the task. For tasks in the BEHAVIOR environment, the dataset also includes accompanying 833 demo videos that showcase the execution of the ground truth action trajectories. Further details 834 regarding the dataset format can be found in Appendix C.4.

835 **How many instances are there in total (of each type, if appropriate)?**

836 We have released a dataset containing task plans and trajectories for two environments: VirtualHome 837 and BEHAVIOR . For VirtualHome , there are 338 task plans/trajectories with a total of 2,960 steps, 838 and 801 goal conditions. For BEHAVIOR , there are 100 task plans/trajectories with a total of 839 1,460 steps, and 673 goal conditions. Additionally, we have annotated transition models for both 840 environments. VirtualHome has 33 annotated transition models with 99 preconditions and 57 effects, 841 while BEHAVIOR has 30 annotated transition models with 84 preconditions and 51 effects.

842 More detailed information can be found in Table 2 of the main paper. Further statistics on the 843 entire dataset are available on the website [https://cs.stanford.edu/~manling1/projects/](https://cs.stanford.edu/~manling1/projects/embodied-eval) 844 embodied-eval. We will continue releasing more data and updating the information on the website in the 845 future.

846 **Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from 847 a larger set?** If the dataset is a sample, then what is the larger set? Is the sample representative of the larger set 848 (e.g., geographic coverage)? If so, please describe how this representativeness was validated/verified. If it is not 849 representative of the larger set, please describe why not (e.g., to cover a more diverse range of instances, because 850 instances were withheld or unavailable).

851 The benchmark is build on top of two simulators VirtualHome [2] and BEHAVIOR -100 [1]. To better evaluate  
852 the decision making ability, we selectively focus on long-horizon tasks with complicated goals in VirtualHome  
853 and BEHAVIOR-100. We select these two simulators due to their length of task plans and the number of goal  
854 conditions, as detailed in Appendix M.1. Detailed statistics of the selected subset is discussed in Appendix M.4,  
855 where we show that the selected subset is very diverse in tasks and environments.

856 **What data does each instance consist of?** Raw data (e.g., unprocessed text or images) or features? In either  
857 case, please provide a description.

858 The benchmark contains task annotations along with an evaluator codebase for executing plans in the simulators <https://github.com/embodied-agent-eval/embodied-agent-eval/>. This codebase  
859 is created by modifying the existing simulator backbones VirtualHome and BEHAVIOR. For example, in the  
860 BEHAVIOR environment, we implemented a symbolic simulator. Detailed information about the evaluator  
861 implementations and the revisions made to the simulators to support the evaluation can be found in Appendix Q.1  
862 for the BEHAVIOR environment and Appendix Q.2 for the VirtualHome environment. Using the provided  
863 simulator, users can visualize the plan execution process and capture screenshots. For tasks in the BEHAVIOR  
864 environment, the dataset also includes demo videos that demonstrate the execution of the ground truth action  
865 trajectories, complementing the annotated ground truth goals and plans.

867 **Is there a label or target associated with each instance?** If so, please provide a description.

868 Each instance is associated with a sequence of labels showing the ground truth action sequence.

869 **Is any information missing from individual instances?** If so, please provide a description, explaining why  
870 this information is missing (e.g. because it was unavailable). This does not include intentionally removed  
871 information but might include, e.g., redacted text.

872 All instances are complete.

873 **Are relationships between individual instances made explicit (e.g., users' movie ratings, social network  
874 links)?** If so, please describe how these relationships are made explicit.

875 Some instances may share the same task name or same goals, but different trajectories. It can be viewed as  
876 alternative decision making trajectories to achieve the goal. It will be indicated by a unique identifier indicating  
877 the scene and the task.

878 **Are there recommended data splits (e.g., training, development/validation, testing)?** If so, please provide a  
879 description of these splits, explaining the rationale behind them.

880 The benchmark presented in this work is specifically designed for zero-shot testing, with the primary goal of  
881 assessing the out-of-the-box long-horizon embodied decision-making capabilities of LLMs.

882 **Are there any errors, sources of noise, or redundancies in the dataset?** If so, please provide a description.

883 The dataset was very carefully manually curated to mitigate any incidence of errors within the goal annotations,  
884 trajectory annotations, and transition model annotations. For VirtualHome , we cross-examine the goal annotations  
885 by executing alternative action trajectories and summarizing the overlapping final states. After that, we  
886 manually cleaned and removed any actions that were deemed unrelated. The annotation process, which utilized  
887 a human-interactive annotation program, is discussed in further detail in Appendix P.2. For the BEHAVIOR  
888 environment, we cross-examine the action sequence annotations by grounding them to multiple demonstration  
889 videos. This process is described in more detail in Appendix P.1.

890 **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets,  
891 other datasets)?** If it links to or relies on external resources, a) are there guarantees that they will exist, and  
892 remain constant, over time; b) are there official archival versions of the complete dataset (i.e., including the  
893 external resources as they existed at the time the dataset was created); c) are there any restrictions (e.g., licenses,  
894 fees) associated with any of the external resources that might apply to a future user? Please provide descriptions  
895 of all external resources and any restrictions associated with them, as well as links or other access points, as  
896 appropriate.

897 The entire dataset will be made publicly available on our project website <https://cs.stanford.edu/~manling1/projects/embodied-eval>. We will also provide a dockerized simulator evaluator tool for  
898 executing plans within the simulator environment. The text data will be released in JSON format and hosted on  
899 our website. Detailed information about the dataset format can be found in Appendix C.4. The benchmark will  
900 be released under the VirtualHome and BEHAVIOR licenses, which allow public use of the simulators for both  
901 research and commercial purposes.

903 **Does the dataset contain data that might be considered confidential (e.g., data that is protected by legal  
904 privilege or by doctor-patient confidentiality, data that includes the content of individuals non-public  
905 communications)?** If so, please provide a description.

906 No.

907 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might**  
908 **otherwise cause anxiety?** If so, please describe why.

909 No.

910 **Does the dataset relate to people?** If not, you may skip the remaining questions in this section.

911 No.

912 **Does the dataset identify any subpopulations (e.g., by age, gender)?** If so, please describe how these  
913 subpopulations are identified and provide a description of their respective distributions within the dataset.

914 No.

915 **Is it possible to identify individuals (i.e., one or more natural persons), either directly or indirectly (i.e., in**  
916 **combination with other data) from the dataset?** If so, please describe how.

917 No.

918 **Does the dataset contain data that might be considered sensitive in any way (e.g., data that reveals racial or**  
919 **ethnic origins, sexual orientations, religious beliefs, political opinions or union memberships, or locations;**  
920 **financial or health data; biometric or genetic data; forms of government identification, such as social**  
921 **security numbers; criminal history)?** If so, please provide a description.

922 No.

923 **Any other comments?**

924 No.

## Collection Process

925 **How was the data associated with each instance acquired?** Was the data directly observable (e.g., raw text,  
926 movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived from other data (e.g.,  
927 part-of-speech tags, model-based guesses for age or language)? If data was reported by subjects or indirectly  
928 inferred/derived from other data, was the data validated/verified? If so, please describe how.

929 We annotate the data based on existing annotations from the VirtualHome and BEHAVIOR environments. To  
930 ensure the accuracy of the annotations, we execute the plans in the simulator and verify that the goals are satisfied.  
931

932

933

934 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, manual**  
935 **human curation, software program, software API)?** How were these mechanisms or procedures validated?

936 As mentioned above, the data is collected based on existing annotations from VirtualHome and BEHAVIOR .  
937 We supplement these annotations with additional manual annotations about goals, trajectories, transition models,  
938 and natural language task instructions, which are described in detail in Appendix P.

939 **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic, probabilistic**  
940 **with specific sampling probabilities)?**

941 To better evaluate the decision making ability, we selectively focus on long-horizon tasks with complicated goals  
942 in VirtualHome and BEHAVIOR-100. We select these two simulators due to their length of task plans and the  
943 number of goal conditions, as detailed in Appendix M.1. Detailed statistics of the selected subset is discussed in  
944 Appendix M.4, where we show that the selected subset is very diverse in tasks and environments.

945 **Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were**  
946 **they compensated (e.g., how much were crowdworkers paid)?**

947 The annotations are purely done by the authors in this paper, which are expert researchers.

948 **Over what timeframe was the data collected? Does this timeframe match the creation timeframe of**  
949 **the data associated with the instances (e.g., recent crawl of old news articles)?** If not, please describe the  
950 timeframe in which the data associated with the instances was created.

951 The newly added annotations are done in 2024. VirtualHome was created in 2018, and BEHAVIOR was created  
952 in 2022.

953 **Were any ethical review processes conducted (e.g., by an institutional review board)?** If so, please provide  
954 a description of these review processes, including the outcomes, as well as a link or other access point to any  
955 supporting documentation.

956 No.

957 **Does the dataset relate to people?** If not, you may skip the remaining questions in this section.

958 No.

959 **Did you collect the data from the individuals in question directly, or obtain it via third parties or other**

960 **sources (e.g., websites)?**

961 N/A.

962 **Were the individuals in question notified about the data collection?** If so, please describe (or show with

963 screenshots or other information) how notice was provided, and provide a link or other access point to, or

964 otherwise reproduce, the exact language of the notification itself.

965 N/A.

966 **Did the individuals in question consent to the collection and use of their data?** If so, please describe (or

967 show with screenshots or other information) how consent was requested and provided, and provide a link or

968 other access point to, or otherwise reproduce, the exact language to which the individuals consented.

969 N/A.

970 **If consent was obtained, were the consenting individuals provided with a mechanism to revoke their**

971 **consent in the future or for certain uses?** If so, please provide a description, as well as a link or other access

972 point to the mechanism (if appropriate).

973 N/A.

974 **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data protection**

975 **impact analysis) been conducted?** If so, please provide a description of this analysis, including the outcomes,

976 as well as a link or other access point to any supporting documentation.

977 N/A.

978 **Any other comments?**

979 No.

## Preprocessing/cleaning/labeling

982 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)? If so, please provide a description. If not, you may skip the remainder of the questions in this section.**

983

984

985 The annotations were curated by humans. We execute the trajectories in the simulators to validate the goals, trajectories, and transition models.

986

987 **Was the raw data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)? If so, please provide a link or other access point to the “raw” data.**

988

989 Human annotations of missing goals, trajectories, and transition models are created from scratch, with the process detailed in Appendix P. The raw demo videos used for annotation reference are publicly available from BEHAVIOR [1].

990

991

992 **Is the software used to preprocess/clean/label the instances available? If so, please provide a link or other access point.**

993

994 The code used to facilitate the dataset annotation process is available at <https://github.com/embodied-agent-eval/embodied-agent-eval/> and is described in detail in Appendix P.

995

996 **Any other comments?**

997 No.

## Uses

1000 **Has the dataset been used for any tasks already?** If so, please provide a description.  
1001 The dataset presented in this work has been used as a benchmark for evaluating the performance of 15 large  
1002 language models (LLMs). More details about these LLMs can be found in Appendix Q.3. Furthermore, the code  
1003 for conducting zero-shot evaluation using our dataset will be made available on our project's GitHub repository.

1004 **Is there a repository that links to any or all papers or systems that use the dataset?** If so, please provide a  
1005 link or other access point.  
1006 It will be made public on our website once more papers start to use our dataset.  
1007 **What (other) tasks could the dataset be used for?** Is there anything about the composition of the dataset or  
1008 the way it was collected and preprocessed/cleaned/labeled that might impact future uses? For example, is there  
1009 anything that a future user might need to know to avoid uses that could result in unfair treatment of individuals or  
1010 groups (e.g., stereotyping, quality of service issues) or other undesirable harms (e.g., financial harms, legal risks)?  
1011 If so, please provide a description. Is there anything a future user could do to mitigate these undesirable harms?  
1012 There may be potential to benchmark vision-language foundation models (VLMs) using our benchmark.  
1013 However, we leave the exploration of these possibilities for future work.  
1014 **Is there anything about the composition of the dataset or the way it was collected and preprocessed/-**  
1015 **cleaned/labeled that might impact future uses?** For example, is there anything that a future user might need  
1016 to know to avoid uses that could result in unfair treatment of individuals or groups (e.g., stereotyping, quality of  
1017 service issues) or other undesirable harms (e.g., financial harms, legal risks)? If so, please provide a description.  
1018 Is there anything a future user could do to mitigate these undesirable harms?  
1019 No.  
1020 **Are there tasks for which the dataset should not be used?** If so, please provide a description.  
1021 No.  
1022 **Any other comments?**  
1023 No.

## Distribution

1024  
1025  
1026 **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization)**  
1027 **on behalf of which the dataset was created?** If so, please provide a description.  
1028 The dataset will be made publicly available and can be used for both research and commercial purposes under  
1029 the VirtualHome and BEHAVIOR licenses.  
1030 **How will the dataset be distributed (e.g., tarball on website, API, GitHub)** Does the dataset have a digital  
1031 object identifier (DOI)?  
1032 The dataset will be distributed as a JSON file that includes a unique identifier for each task goal, along with  
1033 its associated action trajectory, natural language task name, natural language task instructions, symbolic goal  
1034 definition, and the transition models involved in the task. Additionally, we will release the evaluator, which  
1035 enables the automatic calculation of fine-grained systematic metrics. Further details regarding the dataset format  
1036 can be found in Appendix C.4.  
1037 **When will the dataset be distributed?**  
1038 The full dataset will be made available upon the acceptance of the paper before the camera-ready deadline. We  
1039 release it on our website.  
1040 **Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under**  
1041 **applicable terms of use (ToU)?** If so, please describe this license and/or ToU, and provide a link or other access  
1042 point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these  
1043 restrictions.  
1044 The dataset will be publicly released under the VirtualHome and BEHAVIOR licenses, which allows direct  
1045 public use for both research and commercial purposes.  
1046 **Have any third parties imposed IP-based or other restrictions on the data associated with the instances?**  
1047 If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any  
1048 relevant licensing terms, as well as any fees associated with these restrictions.  
1049 No.  
1050 **Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?** If  
1051 so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any  
1052 supporting documentation.  
1053 No.

1054 **Any other comments?**

1055 No.

## Maintenance

1058 **Who will be supporting/hosting/maintaining the dataset?**

1059 The authors of the paper will maintain the dataset, and pointers to the dataset will be hosted on the GitHub repos-  
1060 itory <https://github.com/embodied-agent-eval/embodied-agent-eval/>. The repository  
1061 will also include the code for downloading the dataset and the evaluation tool.

1062 **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

1063 Contact information for the authors is available on our website. We can also be reached through GitHub issues  
1064 or via email for any inquiries or support related to the dataset and evaluation tool.

1065 **Is there an erratum?** If so, please provide a link or other access point.

1066 We will maintain an erratum on the GitHub repository to host any approved corrections or updates suggested by  
1067 the authors or the broader video research community.

1068 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?** If  
1069 so, please describe how often, by whom, and how updates will be communicated to users (e.g., mailing list,  
1070 GitHub)?

1071 Yes, we plan to host an erratum publicly on our Github. We also plan to expand the scale of the annotations, and  
1072 any updates or extensions will be announced on our website and GitHub repository.

1073 **If the dataset relates to people, are there applicable limits on the retention of the data associated with the  
1074 instances (e.g., were individuals in question told that their data would be retained for a fixed period of  
1075 time and then deleted)?** If so, please describe these limits and explain how they will be enforced.

1076 No.

1077 **Will older versions of the dataset continue to be supported/hosted/maintained?** If so, please describe how.  
1078 If not, please describe how its obsolescence will be communicated to users.

1079 N/A. There are no older versions at the current moment. All updates regarding the current version will be  
1080 communicated via our website and Github.

1081 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do  
1082 so?** If so, please provide a description. Will these contributions be validated/verified? If so, please describe  
1083 how. If not, why not? Is there a process for communicating/distributing these contributions to other users? If so,  
1084 please provide a description.

1085 Contributions to the dataset and evaluation tool will be made possible using standard open-source practices.  
1086 Interested parties can submit pull requests to the relevant GitHub repository, which will be reviewed and  
1087 incorporated by the authors.

1088 **Any other comments?**

1089 No.

## C URLs of Data and License

1091 The dataset is uploaded for public download under the CC-BY-4.0 license: <https://cs.stanford.edu/~manling1/projects/embodied-eval>.

### C.1 Hosting and Maintenance Plan

1094 Our dataset is also hosted on <https://github.com/embodied-agent-eval/embodied-agent-eval>/which will provide long-term support for hosting the dataset. Contact  
1095 information for the authors is available on our website. We can also be reached through GitHub issues or via  
1096 email for any inquiries or support related to the dataset and evaluation tool. We will maintain an erratum on  
1097 the GitHub issues to host any approved corrections or updates suggested by the authors or the broader video  
1098 research community.

### C.2 Long-term Preservation and DOI

1101 We provide a persistent dereferenceable identifier DOI: [https://datadryad.org/stash/share/bLEj8IHqyKc9c\\_Ff0M8tTSZMTxOLxIup5zUWn9yq\\_j8](https://datadryad.org/stash/share/bLEj8IHqyKc9c_Ff0M8tTSZMTxOLxIup5zUWn9yq_j8).

1103 **C.3 URL of Croissant Metadata Record**

1104 We provide structured metadata (schema.org standards) in <https://github.com/embodied-agent-eval/embodied-agent-eval/blob/main/dataset/metadata.json>.

1106 **C.4 Dataset Format**

1107 Each instance in the dataset represents a task goal. Specifically, each task contains the following data: (1) 1108 natural language task name, (2) natural language task instruction, (3) symbolic goal definition (including its LTL 1109 form), (4) symbolic action trajectory, (5) the transition models involved in the task. For tasks in the BEHAVIOR 1110 environment, the dataset also includes accompanying VR human demonstration videos that showcase the 1111 execution of the ground truth action trajectories.

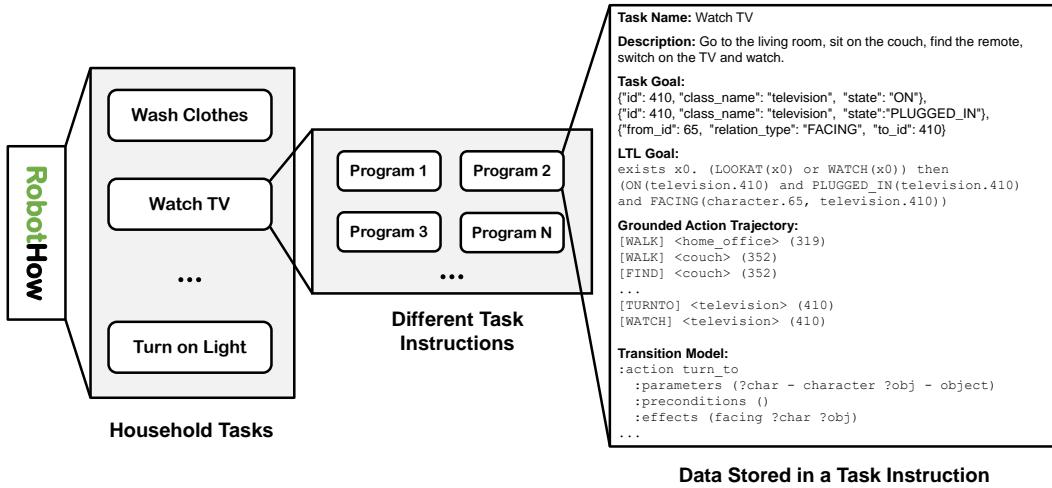


Figure 6: VirtualHome dataset structure example.

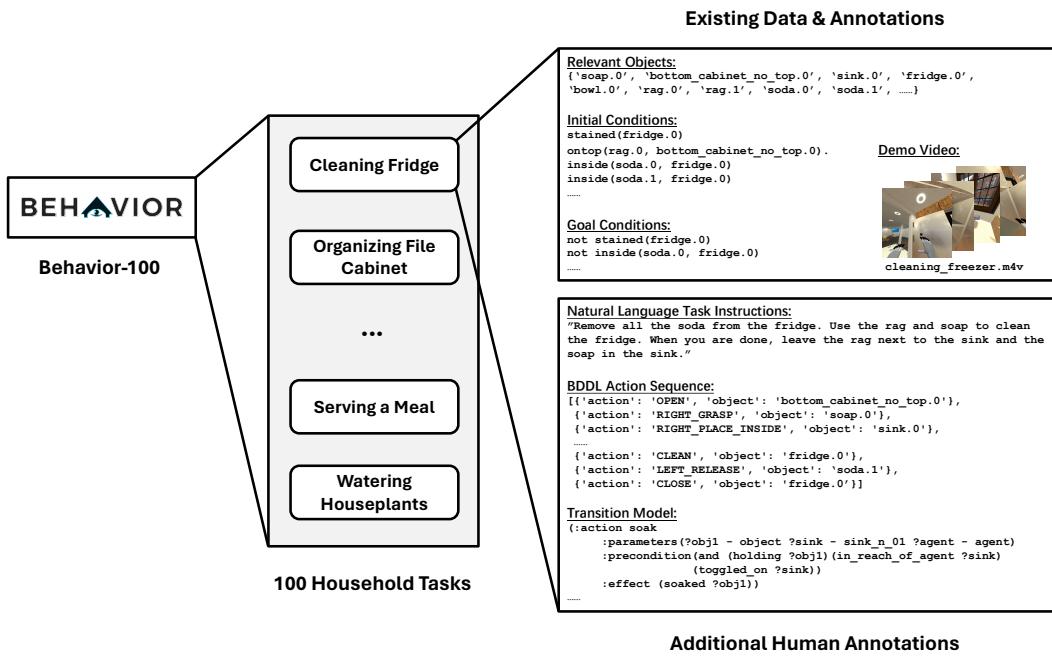


Figure 7: BEHAVIOR dataset structure example.

1112 Please find our JSON data format in this link: <https://github.com/embodied-agent-eval/embodied-agent-eval/tree/main/dataset>:

### Data Format Examples for VirtualHome

```

# Below is an example of VirtualHome
"1057_1": {
    "task_name": "Watch TV",
    "natural_language_description": "Go to the living room, sit on the couch, find the
        remote, switch on the TV and watch",
    "vh_goal": {
        "actions": [
            "LOOKAT|WATCH"
        ],
        "goal": [
            {
                "id": 410,
                "class_name": "television",
                "state": "ON"
            },
            {
                "id": 410,
                "class_name": "television",
                "state": "PLUGGED_IN"
            },
            {
                "from_id": 65,
                "relation_type": "FACING",
                "to_id": 410
            }
        ]
    },
    "tl_goal": "(exists x0. ((LOOKAT(x0) or WATCH(x0))) then (ON(television.410) and
        PLUGGED_IN(television.410) and FACING(character.65, television.410)))",
    "action_trajectory": [
        "[WALK] <home_office> (319)",
        "[WALK] <couch> (352)",
        "[FIND] <couch> (352)",
        "[SIT] <couch> (352)",
        "[FIND] <remote_control> (1000)",
        "[FIND] <television> (410)",
        "[SWITCHON] <television> (410)",
        "[TURNTO] <television> (410)",
        "[WATCH] <television> (410)"
    ],
    "transition_model": <pddl_definition>
}

```

1114

Figure 8: We release the task annotations in JSON format. Here is the data instance example for VirtualHome .  
1115  
1116

### Data Format Examples for BEHAVIOR

```

# Below is an example of BEHAVIOR
"cleaning_high_chair_0_Wainscott_0_int_0_2021-06-05_18-03-15": {
    "task_name": "cleaning_high_chair",
    "natural_description": "Clean the high chair.",
    "raw_bddl_goal": "(define (problem cleaning_high_chair_0)\n    (:domain igibson)\n\n    (:objects\n        highchair.n.01_1 - highchair.n.01\n        piece_of_cloth.n.01_1 - piece_of_cloth.n.01\n        cabinet.n.01_1 - cabinet.n.01\n        ...",
    "tl_goal": "not dusty(highchair.0)",
    "action_trajectory": [
        {
            "action": "OPEN",
            "object": "bottom_cabinet_no_top_80"
        },
        {
            "action": "RIGHT_GRASP",
            "object": "paper_towel_0"
        },
        {
            "action": "LEFT_GRASP",
            "object": "highchair_0"
        },
        {
            "action": "CLEAN",
            "object": "highchair_0"
        }
    ],
    "transition_model": <pddl_definition>,
    "demo": <demo_link>
}

```

1117

1118 Figure 9: We release the task annotations in JSON format. Here is the data instance example for BEHAVIOR .  
1119

## 1120 **C.5 Author Statement**

1121 The authors bear all responsibility in case of violation of rights. All dataset annotations were collected by the  
1122 authors and we are releasing the dataset under CC-BY-4.0.

## 1123 **D URLs of Code and Re-productivity**

1124 We have made the codebase publicly available on GitHub (<https://github.com/embodied-agent-eval/embodied-agent-eval/>), along with detailed instructions. In this  
1125 section, we provide the URLs for the VirtualHome evaluator, BEHAVIOR evaluator, and LLM implementations.  
1126 Further details regarding these implementations can be found in Appendix Q.

### 1128 **D.1 Code for VirtualHome Evaluator and Computation Resources**

1129 We provide the code for the VirtualHome evaluator at <https://github.com/embodied-agent-eval/embodied-agent-eval/tree/main/virtualhome>, along with detailed usage instructions. This code  
1130 facilitates goal interpretation, subgoal decomposition, action sequencing, and transition modeling.  
1131

1132 To use our evaluator, ensure that Docker is installed in your environment. For annotating the RobotHow dataset,  
1133 we also offer code for human interactive annotation of goals within VirtualHome , with usage instructions  
1134 available in the repository.

1135 The experiment of subgoal decomposition is run on a Windows 11 system with an AMD Ryzen 7 5800H  
1136 processor and 16GB of memory. For goal interpretation, action sequencing and transition modeling, the  
1137 experiments are run on a Linux system with 8 Intel(R) Xeon(R) Platinum 8481C CPU @ 2.70GHz and 32GB of  
1138 memory.

### 1139 **D.2 Code for BEHAVIOR Evaluator and Computation Resources**

1140 The code for the BEHAVIOR evaluator is available on GitHub at <https://github.com/embodied-agent-eval/embodied-agent-eval/tree/main/behavior>. We provide easy-to-  
1141 reproduce scripts for prompt generation and evaluation of each module: *goal interpretation, action sequencing,*  
1142 *subgoal decomposition, and transition modeling*. The BEHAVIOR related experiments were conducted on a  
1143 Windows 11 system equipped with an AMD Ryzen 7 7800X3D processor and 32GB of RAM.  
1144

### 1145 **D.3 Code for LLMs Implementations and Computation Resources**

1146 We provide convenient and easy-to-reproduce batch LLM inference by directly integrating our evaluation  
1147 pipeline into the the HELM [67] code base. In order to allow quick and easy access for all users with different  
1148 development environment and GPU setups, we provide the option to use Together AI <sup>†</sup> APIs for large open-source  
1149 model inference. Users can easily setup their inference environment by following the instructions in <https://github.com/embodied-agent-eval/embodied-agent-eval/tree/main/helm>.  
1150

1151 To achieve standardization and consistency across all models, we used the same set of decoding parameters  
1152 for all LLMs evaluated. Specifically, we used temperature zero for all models because we wanted to use the  
1153 arg max under the model’s distribution. Also, because several of the models only support temperature-based  
1154 sampling, as opposed to other sampling methods, we restricted ourselves to temperature-scaling during the  
1155 sample process. For a given prompt, since the model’s completion involves sampling, there is randomness  
1156 involved in determining the specific completion decoded for each instance. For our scenarios, since we are  
1157 decoding the arg max through low temperature decoding, this is not a significant factor.

1158 To perform a single run of all models on our benchmark, a total of 180 runs would be required. This involves  
1159 evaluating each specific model on each specific simulator ability module. The total cost of this process amounts  
1160 to approximately 126,900,000 tokens and 50,760 queries across all models, which sums to a compute cost of  
1161 \$879.24. For open-source models, the costs are calculated based on the pricing of the Together.ai API.

1162 Due to space limitations, we use shortened model names for common commercial and open-source LLMs  
1163 evaluated in the main paper. However, shortened model names do not specify the precise release version  
1164 and hosting organization of the models being addressed. To support maximum reproducibility, we list detail  
1165 information for each model in Table 9.

---

<sup>†</sup><https://www.together.ai/>

Table 9: Model Cards for All Evaluated Large Language Models

Model Name	Creator	Complete Model ID	Release	Hosting
Claude-3 Haiku	Anthropic	claude-3-haiku-20240307	03/07/24	Anthropic
Claude-3 Sonnet	Anthropic	claude-3-sonnet-20240229	02/29/24	Anthropic
Claude-3 Opus	Anthropic	claude-3-opus-20240229	02/29/24	Anthropic
Cohere Command R	Cohere	command-r	03/11/24	Cohere
Cohere Command R+	Cohere	command-r-plus	04/04/24	Cohere
Gemini 1.0 Pro	Google	gemini-pro	12/13/23	GCP Vertex
Gemini 1.5 Flash	Google	gemini-1.5-flash-preview-0514	05/14/24	GCP Vertex
Gemini 1.5 Pro	Google	gemini-1.5-pro-preview-0409	04/09/24	GCP Vertex
GPT-3.5-turbo	OpenAI	gpt-3.5-turbo-0125	01/25/24	OpenAI
GPT-4-turbo	OpenAI	gpt-4-turbo-2024-04-09	04/09/24	OpenAI
GPT-4o	OpenAI	gpt-4o-2024-05-13	05/13/24	OpenAI
Llama3 8B Instruct	Meta	meta-llama-3-8b-instruct	04/18/24	TogetherAI
Llama3 70B Instruct	Meta	meta-llama-3-70b-instruct	04/18/24	TogetherAI
Mistral Large	MistralAI	mistral-large-2402	02/26/24	MistralAI
Mixtral 8x22B MoE	MistralAI	mixtral-8x22b-instruct-v0.1	04/17/24	TogetherAI

## 1166 E Impact, Limitations and Future Directions

### 1167 E.1 Broader Impact

1168 The proposed EMBODIED AGENT INTERFACE and the comprehensive evaluation of Large Language Models  
 1169 (LLMs) for embodied decision-making have significant implications for the development and deployment of  
 1170 intelligent agents in various domains, including robotics, autonomous systems, and human-robot interaction.

1171 Recent advancements in LLMs have revolutionized the field of artificial intelligence, particularly in the domains  
 1172 of digital and embodied agents. LLMs have demonstrated remarkable capabilities in natural language under-  
 1173 standing, generation, and reasoning, enabling the development of sophisticated AI systems that can interact with  
 1174 humans and the environment in more natural and intuitive ways. The potential of LLMs extends beyond purely  
 1175 digital interactions, as embodied agents that can perceive, reason, and act within physical environments have  
 1176 also begun to benefit from the integration of LLMs. By combining the language understanding and generation  
 1177 capabilities of LLMs with the perceptual and motor skills of embodied agents, researchers aim to create more  
 1178 versatile and intelligent robots that can assist humans in various tasks.

1179 However, despite the progress made in applying LLMs to embodied agents, significant challenges remain.  
 1180 Embodied agents require the ability to ground natural language instructions in the physical world, reason  
 1181 about object properties and relationships, and plan and execute actions in dynamic environments. While LLMs  
 1182 have shown promise in these areas, they exhibit inconsistencies in their interactions with the physical world,  
 1183 occasionally making correct decisions but frequently producing incorrect plans without reliable ways to control  
 1184 their performance. This inconsistency leads to doubts regarding their suitability and readiness for robotic tasks.  
 1185 Current evaluations of LLMs in embodied decision-making tasks often entangle various capabilities, making it  
 1186 difficult to understand why LLMs sometimes work and sometimes fail. Existing benchmarks oversimplify the  
 1187 problem by focusing only on high-level semantic failures, assuming that low-level physics can be automatically  
 1188 fulfilled. This assumption results in unrealistic plans, such as a robot heating food in a microwave without first  
 1189 ensuring the door is closed or the food is properly placed inside, which ignores crucial physical preconditions  
 1190 and state changes.

1191 The standardization of goal specifications, modules, and interfaces through the EMBODIED AGENT INTERFACE  
 1192 framework enables a more systematic and rigorous evaluation of LLMs' capabilities in embodied decision-  
 1193 making tasks. This standardization facilitates the comparison of different approaches and architectures, promoting  
 1194 the development of more advanced and reliable LLM-based embodied agents. By identifying the strengths and  
 1195 weaknesses of current LLMs through fine-grained metrics and error analysis, researchers and practitioners can  
 1196 focus their efforts on addressing specific limitations and improving the overall performance of these agents.

1197 Moreover, the insights gained from the evaluation of LLMs in embodied decision-making tasks can inform the  
 1198 design of future LLMs and training strategies. The identified challenges, such as the difficulty in grounding  
 1199 natural language instructions to environment-specific objects and states, the lack of reasoning abilities in handling  
 1200 preconditions and post-effects of actions, and the reporting bias in goal interpretation, highlight the need for  
 1201 more targeted training approaches. These findings can guide the development of LLMs that are better suited for  
 1202 embodied decision-making, potentially leading to more effective and efficient agents in real-world applications.  
 1203 The fine-grained evaluation results and the breakdown of LLM abilities into detailed tests can pinpoint the  
 1204 strengths and weaknesses in LLM interactions with the physical world, offering a more robotics-centered  
 1205 evaluation to uncover how much LLMs understand about the physical world and what knowledge can be  
 1206 improved or added for robotics.

1207 The EMBODIED AGENT INTERFACE framework and the fine-grained evaluation results can also contribute to the  
1208 development of safer and more trustworthy embodied agents. By providing a comprehensive understanding of  
1209 LLMs' limitations and failure modes in embodied decision-making, this work can help in designing safeguards  
1210 and fallback mechanisms to mitigate potential risks associated with the deployment of these agents in real-world  
1211 scenarios. The fine-grained error analysis can inform the development of monitoring and intervention strategies  
1212 to ensure the safe and reliable operation of LLM-based embodied agents.

1213 Furthermore, the EMBODIED AGENT INTERFACE framework and the evaluation methodology can be extended  
1214 to other domains beyond household tasks, such as industrial automation, healthcare robotics, and autonomous  
1215 vehicles. The standardization of goal specifications, modules, and interfaces can facilitate the application of  
1216 LLMs in these domains, enabling the development of more intelligent and adaptive agents that can understand  
1217 and execute complex tasks based on natural language instructions.

1218 However, it is essential to consider the potential negative impacts and ethical implications of deploying LLM-  
1219 based embodied agents. The evaluation results highlight the need for careful consideration of the limitations  
1220 and biases of these agents, particularly in safety-critical applications. It is crucial to establish guidelines and  
1221 best practices for the responsible development and deployment of LLM-based embodied agents, ensuring  
1222 transparency, accountability, and alignment with human values.

1223 In conclusion, the EMBODIED AGENT INTERFACE framework and the comprehensive evaluation of LLMs for  
1224 embodied decision-making have the potential to advance the field of intelligent agents and enable the development  
1225 of more capable, reliable, and trustworthy embodied agents. By providing a standardized approach for evaluation  
1226 and identifying key challenges and opportunities, this work can guide future research and development efforts in  
1227 this area, ultimately leading to the realization of intelligent agents that can effectively understand and execute  
1228 complex tasks in real-world environments. As LLMs continue to evolve and be integrated into embodied agents,  
1229 it is essential to address the challenges and ethical considerations to ensure their responsible and beneficial  
1230 deployment in various domains.

## 1231 **E.2 Limitations**

1232 While we currently evaluate the capabilities and limitations of LLMs in embodied decision-making tasks, our  
1233 approach has limitations as we abstract input environments using relational graphs of objects, which may not  
1234 adequately represent the rich multimodal information, low-level physical properties, and dynamics of real-world  
1235 environments. It does not cover the challenge of grounding natural language instructions in perceptual data and  
1236 executing actions that require fine-grained control and precise manipulation. Moreover, our current evaluation  
1237 primarily focuses on symbolic reasoning and decision-making, without directly incorporating sensory inputs or  
1238 actuation outputs. While this approach allows us to isolate and assess specific cognitive capabilities, it neglects  
1239 the crucial integration of perception and action, which is fundamental to embodied agents operating in physical  
1240 environments. Realistic embodied decision-making requires seamless interplay between language understanding,  
1241 sensory processing, and motor control, which our current framework does not fully capture.

## 1242 **E.3 Potential Negative Social Impact**

1243 The development of embodied agents powered by large language models has significant potential social impacts  
1244 that must be carefully considered. On one hand, these intelligent agents could revolutionize numerous domains  
1245 and provide immense societal benefits. In healthcare, they could assist in elderly care, rehabilitation, and  
1246 medical procedures. In education, they could serve as personalized companions and learning companions. In  
1247 manufacturing and logistics, they could streamline operations, improve efficiency, and enhance workplace safety.

1248 However, the deployment of such capable agents also raises concerns over privacy, security, and ethical  
1249 implications. There are risks of these systems being exploited for malicious purposes or exhibiting harmful  
1250 biases learned from training data. Additionally, their widespread adoption could disrupt certain job markets and  
1251 exacerbate economic inequalities if not managed responsibly. As we strive to unlock the full capabilities of  
1252 large language models for embodied decision-making, it is crucial to prioritize the development of robust safety  
1253 measures, ethical frameworks, and regulatory guidelines. Engaging diverse stakeholders, including policymakers,  
1254 domain experts, and the general public, will be vital in navigating the intricate societal impacts and ensuring  
1255 these powerful technologies are harnessed for the greater benefit of humanity.

## 1256 **E.4 Potential Future Directions**

1257 We outline the future directions to represent a roadmap for advancing embodied decision making with large  
1258 foundation models.

- 1259 • **Multimodal Grounding and Integration:** Extend the evaluation to incorporate multimodal inputs  
1260 such as vision, audio, by integrating with Vision-Language Models (VLMs) and other multimodal  
1261 models. This includes assessing LLMs' ability to ground language in rich sensory contexts, generate  
1262 low-level control commands, and reason about visual information like object keypoints, grasp poses,  
1263 and scene configurations. Additionally, explore the use of symbolic scene graph representations as  
1264 input to LLMs, enabling more effective reasoning about object relationships, spatial configurations,  
1265 and visual attributes.

- **Vision and Low-Level Control:** Evaluate LLMs in scenarios that require tight integration with vision systems and low-level control modules, such as tabletop object manipulation tasks. This could involve predicting object keypoints, generating pick-and-place grasp poses, or generating goal configurations for object arrangements.
- **Episodic Memory and Scene Graph Memory:** Investigate the integration of memory systems into LLMs, including episodic memory for storing and retrieving past experiences, state memory for maintaining internal environment representations, and scene graph memory structures. Evaluating memory-augmented LLMs can provide insights into leveraging past knowledge and rich environment representations for improved decision-making.
- **Physical and Geometric Reasoning:** Incorporate tasks that require physical and geometric reasoning capabilities, such as spatial planning, object manipulation, collision avoidance, stability prediction, and reasoning about object affordances, sizes, shapes, and attachments. Evaluate LLMs' potential for precise physical interactions and reasoning in real-world applications.
- **Navigation and Exploration:** Include navigation tasks and scenarios that require LLMs to understand navigational instructions, plan efficient routes, adapt to dynamic environments, and explore unknown object states, physical properties, and environmental conditions crucial for robust decision-making.
- **Forward, Backward, and Counterfactual Prediction:** Assess LLMs' ability to perform forward prediction (action → state change), backward prediction (goal state → necessary objects and subgoals), and counterfactual reasoning (hypothetical "what-if" scenarios) to evaluate their capability in multi-step planning and reasoning about the effects of actions or action sequences.
- **Dataset, Simulation, and Benchmark Diversity:** Develop new datasets, simulation environments, and integrate with existing robotics benchmarks to capture the complexity and diversity of real-world scenarios. This includes creating environments with richer physics simulations, more diverse object interactions, and more challenging task setups, as well as integrating with benchmarks focused on task and motion planning, manipulation, and navigation.
- **Fine-Tuning Decision Making Models:** Explore fine-tuning LLMs using the fine-grained error analysis and feedback from the evaluation framework, enabling them to learn from their mistakes and improve decision-making capabilities. Additionally, develop specialized models such as a decision making GPT, trained on trajectories of embodied decision-making tasks, to generalize to unseen trajectories and environments.

By addressing these directions, researchers and practitioners can overcome the current limitations, bridge the gap between language models and embodied agents, and unlock the full potential of LLMs in real-world, dynamic environments. Ultimately, these efforts will pave the way for more capable, adaptable, and trustworthy embodied agents that can seamlessly understand natural language instructions, reason about physical properties and constraints, and execute complex tasks with precision and robustness.

## 1301 F Embodied Agent Interface Details

1302 We will introduce the additional details about the EMBODIED AGENT INTERFACE (EAGENT) in this section,  
1303 including the motivation of the current design and its relationship with Markov Decision Process.

### 1304 F.1 Input and Output Details

1305 As shown in Figure 10, the overall input of the interface consists of three main parts: (1) the task name and  
1306 instruction, (2) the agent instructions, including in-context examples, and (3) the environment representation,  
1307 which includes objects, their states, and relations. The detailed prompt templates are provided in Appendix L.

1308 The input and output for each ability module differ, as illustrated in Figure 11. The mathematical formulation  
1309 has been detailed in Section 2 of the main paper. **Goal Interpretation (ability module 1)** aims to ground the  
1310 natural language instruction to the environment representations of objects, states, relations, and actions. For  
1311 example, in Figure 11, the task instruction "*Use the rag to clean the trays, the bowl, and the refrigerator. When*  
1312 *you are done, leave the rag next to the sink...*" can be grounded to specific objects with IDs, such as *fridge* (ID:  
1313 97), *tray* (ID: 1), *bowl* (ID: 1), *rag* (ID: 0), and *sink* (ID: 82). Note that a simple natural language description can  
1314 be grounded into a set of multiple goal conditions (object state and relation).

1315 The **Subgoal Decomposition (ability module 2)** generates a sequence of states, where each state can be  
1316 a set of objects and their states. Here, we highlight the important states, such as the transitions between a  
1317 sequence of *next\_to(rag.0, sink.82)*, *toggled\_on(sink.82)*, *soaked(rag.0)*, *toggled\_off(sink.82)*, *open(fridge.97)*,  
1318 *not\_stained(fridge.97)*. To achieve these state transitions, we can use a high-level planner such as BFS to search  
1319 for the **Action Sequences (ability module 3)** that achieve these state transitions. We obtain the following action  
1320 sequence: *RIGHT\_GRASP(rag.0)*, *RIGHT\_PLACE\_NEXTTO(sink.82)*, *TOGGLE\_ON(sink.82)*, *SOAK(rag.0)*,  
1321 *TOGGLE\_OFF(sink.82)*, *OPEN(fridge.97)*, *CLEAN(fridge.97)*. Note that multiple actions may be required to  
1322 achieve a single one-step state transition. For example, to perform the state transition *next\_to(rag.0, sink.82)* →

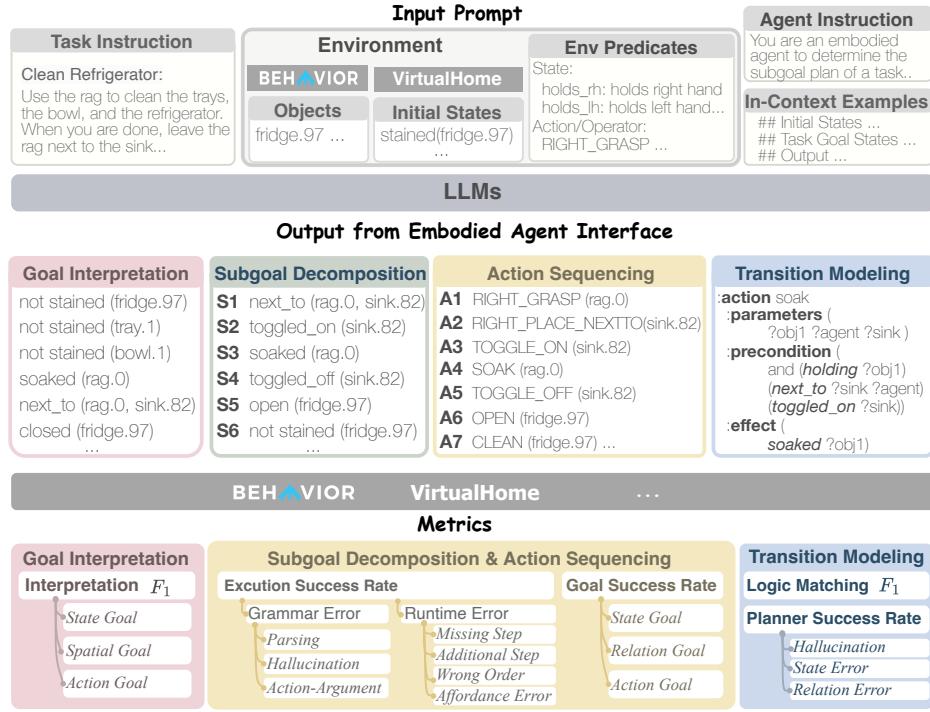


Figure 10: Example input and output for the ability modules.

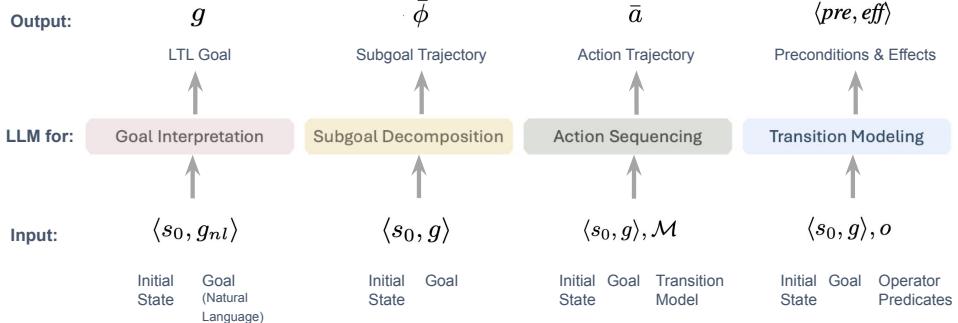


Figure 11: The input and output formulation of four ability modules.

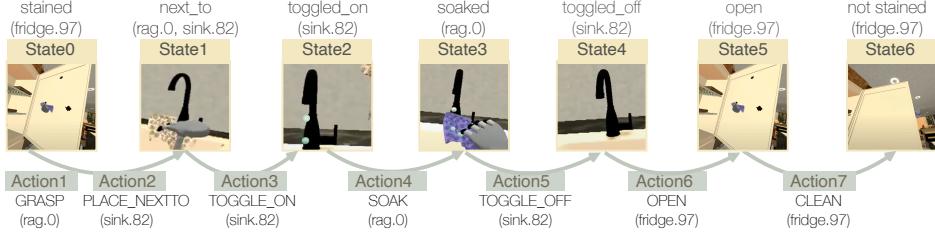
1323 *toggled\_on(sink.82)*, we need two actions *RIGHT\_GRASP(rag.0)*, *RIGHT\_PLACE\_NEXTTO(sink.82)*. We  
1324 show a successful excution of this piece of action sequence in Figure 12.

1325 **Transition Modeling (ability module 4)** is different from the previous modules. It serves as the low-level  
1326 controller to guide the simulator in performing state transitions from preconditions to post-effects. In Figure 11,  
1327 the input is the operator name “*soak*”, and the preconditions are three states: “*holding* (?obj1)”, “*next\_to* (?sink  
1328 ?agent)”, and “*toggled\_on* (?sink)”. The post effect after executing *SOAK* is “*soaked* (?obj1)”.

## 1329 F.2 Grounding to Markov Decision Process

1330 To support a wide range of tasks in various environments, we design the EMBODIED AGENT INTERFACE based  
1331 on the Markov Decision Process (MDP) [68], a fundamental mathematical framework for robot learning to  
1332 formalize sequential decision making in embodied agents [69]. This allows us to create a structured approach to  
1333 benchmark the robot’s decision-making process.

1334 An embodied agent takes natural language instructions from humans and achieves the specified goals through a  
1335 sequence of physical state transitions. It is essentially a decision-making process to determine the actions based  
1336 on the goal and the current state in the embodied environment. As a result, we formulate the MDP process as  
1337 below to input natural language instructions and interact with the environment to achieve the specified goals.



**Example Task** Clean Refrigerator: use the rag to clean the refrigerator and ...

Figure 12: An example of successful execution in BEHAVIOR .

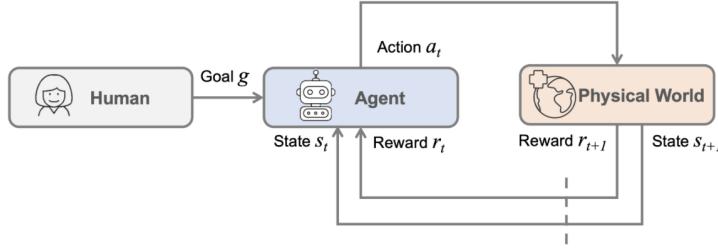


Figure 13: Embodied Decision Making is a Markov Decision Process.

1338 **MDP Formulation for Embodied Agents.** As shown in Figure 13, the Markov Decision Process for an  
 1339 embodied agent can be defined by a tuple  $\langle \mathcal{U}, \mathcal{S}, \mathcal{A}, \mathcal{M}, \mathcal{R}, g \rangle$ , where:

1340  $\mathcal{U}$  is the universe of objects in the environment, which are the fundamental entities that the agent interacts with.  
 1341  $\mathcal{S}$  is the state space, where each state  $s \in \mathcal{S}$  is represented as a tuple  $\langle \mathcal{U}, \mathcal{F} \rangle$ .  $\mathcal{F}$  is a set of relational Boolean  
 1342 features that capture the properties and relations among objects in the environment.  $\mathcal{A}$  is the action space, which  
 1343 represents the set of actions the embodied agent can execute. Actions are represented as tuples  $\langle name, args \rangle$ ,  
 1344 where *name* is the action name and *args* are the object arguments the action operates on.  $\mathcal{M} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$  is  
 1345 the environmental transition model, which specifies the next state  $s_{t+1}$  given the current state  $s_t$  and action  $a$ .  
 1346  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \times g \rightarrow \mathbb{R}$  is the reward function. It depends on the current state, action, and the goal specification  
 1347  $g$ . For a state  $s$ , action  $a$ , and goal  $g$ ,  $\mathcal{R}(s, a, g) = 1$  if  $eval(g, s) = 1$  (i.e., the goal is satisfied in state  $s$ ),  
 1348 and  $\mathcal{R}(s, a, g) = 0$  otherwise. Here,  $eval : g \times \mathcal{S} \rightarrow \{0, 1\}$  determines whether a state satisfies the goal  
 1349 specification.  $g$  is the goal specification. The goal should be grounded in terms of desired final states of objects  
 1350 and their interactions (relations and executed actions), capturing the intended outcome of the agent's actions.  
 1351 The input of goal can be natural language, such as “cleaning the refrigerator” or “polishing furniture”. We  
 1352 denote natural language goal specification as  $g_{nl}$ .

1353 **Grounding Our Evaluation Protocol to the Fundamental Modules of MDP.** The embodied agent receives a  
 1354 natural language goal specification  $g_{nl}$ , translates it to the environment objects and their states, relations, and  
 1355 actions as a goal specification  $g$ , and aims to achieve it through a sequence of state transitions. To abstract the  
 1356 embodied environment, we design the representation to contain *Object*, *State*, *Action*, and, based on that, *Goal*  
 1357 (as final states) and *Trajectory* (as temporally dependent sequences of actions/states). Our interface is built upon  
 1358 a LTL REPRESENTATION layer based on Linear Temporal Logic (LTL), which serves as a unified, expressive  
 1359 interface to communicate with robots in different environments (e.g., different simulators such as BEHAVIOR  
 1360 and VirtualHome).

1361 At each step, the agent observes the current state  $s \in \mathcal{S}$ , selects an action  $a \in \mathcal{A}$  based on its policy  $\pi : \mathcal{S} \times g \rightarrow$   
 1362  $\mathcal{A}$ , and receives a reward  $\mathcal{R}(s, a, g)$ . The environment transitions to the next state according to  $\mathcal{M}(s, a)$ . As  
 1363 shown in Figure 14, according to MDP, it essentially focuses on four abilities:

- 1364 • Input of Goals, which corresponds to **Goal Interpretation (ability module 1)**, translating the natural  
 1365 language goal to environment objects and their relations and actions.
- 1366 • Output of Trajectories, where the output can be a sequence of actions or a sequence of states, which  
 1367 can be regarded as **Action Sequencing (ability module 2)** and **Subgoal Decomposition (ability**  
 1368 **module 3)**.
- 1369 • The core part of the **Transition Model (ability module 4)** Learning, which is covered by the Transition  
 1370 Modeling (ability module 4).

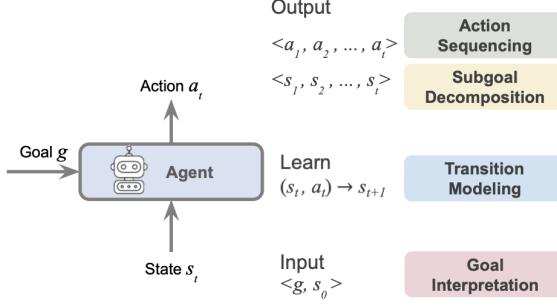


Figure 14: Our four ability modules are fundamental modules of the MDP process.

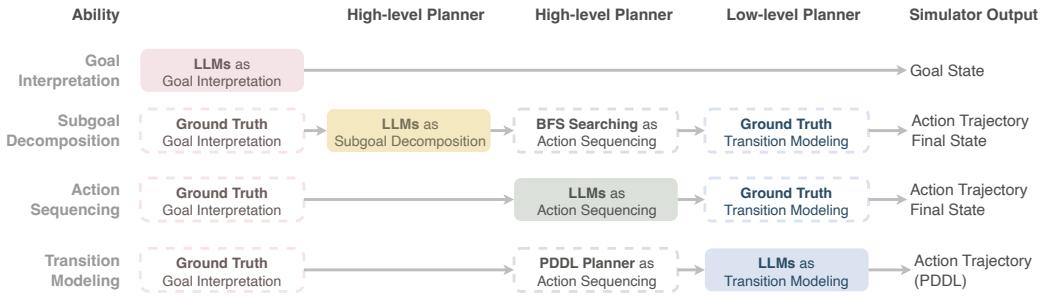


Figure 15: The overview of evaluation pipeline for four abilities. To provide a detailed evaluation of each ability module, we isolate a single module to be handled by the LLMs while using existing data or tools for the other modules. We also conduct a pipeline-based vs modularized analysis, detailed in the Appendix J.

- The goal-evaluation function *eval* and reward model can be reflected in the detailed, fine-grained evaluation metrics we provide.
- 1371     In this way, the EMBODIED AGENT INTERFACE has a comprehensive coverage of the fundamental abilities and  
 1372     can provide a systematic evaluation of the foundational MDP process.

### 1375 F.3 The Relationship between Ability Modules

- 1376 To identify the weaknesses and areas for improvement in LLMs for embodied decision making, we need to  
 1377 evaluate each ability module individually and focus on detailed, fine-grained tasks. Rather than simply knowing  
 1378 that the final success rate is still insufficient, we aim to understand which abilities are already well-developed  
 1379 and how we can effectively integrate different ability modules to enhance overall performance. This includes  
 1380 exploring the integration between LLMs and external tools, as well as LLMs across different modules, enabling  
 1381 us to guide embodied agents to use LLMs more selectively and effectively.
- 1382 To achieve this, as shown in Figure 15, we design an evaluation protocol that isolates a single module to be  
 1383 handled by the LLM while using existing data or tools to serve as the other modules. This approach shifts  
 1384 the focus from an end-to-end evaluation to an accurate assessment of each individual component. By doing  
 1385 so, we can probe the LLM’s capabilities and limitations within each specific ability in detail, gaining a more  
 1386 nuanced understanding of its performance. This fine-grained evaluation allows us to identify the strengths and  
 1387 weaknesses of LLMs in each ability module, guiding future research efforts to address the identified challenges  
 1388 and improve the integration of LLMs in embodied decision-making tasks.
- 1389 The Subgoal Decomposition and Action Sequencing modules are similar in that they both involve trajectory  
 1390 output and evaluate the ordering of decision making. However, the fundamental distinction between them lies  
 1391 in the nature of their outputs. Action sequencing produces imperative actions, while subgoal decomposition  
 1392 generates declarative states, as illustrated in Figure 12.
- 1393 Transition modeling can be considered as the low-level controller that governs the state transitions when  
 1394 executing an action. The hallmark for transition modeling is the ability to search a path to navigate from initial  
 1395 predicates to goal predicates using existing actions. Defining preconditions and post effects for each action  
 1396 enables this search and backtracking.

## 1397 G Fine-Grained Metrics and Automatic Error Detection

1398 To evaluate each ability in the simulator, we design the evaluation pipeline of each ability and detailed in this  
 1399 section.

### 1400 G.1 Goal Interpretation: State Goal, Relation Goal and Action Goal

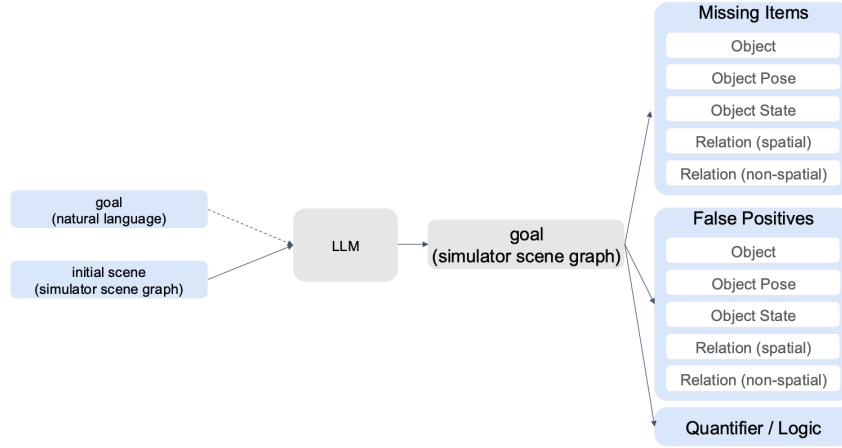


Figure 16: Evaluation pipeline of goal interpretation. We evaluate the LLM’s output for the Goal Interpretation in three key dimensions: single-object states, object-object relations, and agent-action goals. For each dimension, we calculate precision, recall, and F1 score to measure the LLM’s false positive predictions, likelihood of missing goals, and overall capability, respectively.

1401 For the Goal Interpretation, the LLMs receive a natural language description of the overall goal, the initial world  
 1402 state  $S_0$  (including relevant objects and their initial states), and the complete list of all possible actions and  
 1403 object states recognized by the simulator. Based on these inputs, the LLMs are expected to output a symbolic  
 1404 representation of the overall goal, denoted as  $g$ . This task is designed to assess the LLMs’ ability to act as a  
 1405 translation layer between a human user (who only interacts with the system using natural language) and the  
 1406 embodied agent (which only understands symbolic goals composed of simulator-recognizable states and relevant  
 1407 objects in the scene).

1408 To evaluate the LLMs’ output for the Goal Interpretation, we first filter for grammatically correct predicted goals  
 1409 while keeping track of structurally incoherent predictions and object/state hallucinations. Then, we evaluate the  
 1410 remaining grammatically correct predicted goals against ground truth goals in three key dimensions:

- 1411 • **State Goal (or Node Goal)** with a focus on single-object states, e.g., *facing(cat)*, *clean(cup)*,  
 1412 *switched\_on(television)*.
- 1413 • **Relation Goal (or Edge Goal)** with a focus on object-object relations, e.g., *next\_to(character, cat)*,  
 1414 *on(cup, table)*, *on(character, sofa)*.
- 1415 • **Action Goal** with a focus on agent-action goals, e.g., *touch(cat)*, *touch(remote)*. Simulators such as  
 1416 VirtualHome contain action goals for some short tasks that have key actions but no post-effects to  
 1417 validate in the goal, such as the task of “*pat cat*”. In contrast, other simulators like BEHAVIOR-100  
 1418 do not include such action goals.

1419 For each of the three goals, we calculate the following metrics:

- 1420 • *Precision*: Measures the LLMs’ false positive predictions, indicating the proportion of predicted goals  
 1421 that are correct.
- 1422 • *Recall*: Measures the LLMs’ likelihood to miss certain goals, indicating the proportion of ground truth  
 1423 goals that are correctly predicted.
- 1424 • *F<sub>1</sub>* Score: A joint measure that combines precision and recall, representing the overall capability of  
 1425 the LLMs in each dimension.

1426 These metrics provide a comprehensive evaluation of the LLMs’ performance in interpreting natural language  
 1427 goals and translating them into symbolic representations that the embodied agent can understand and execute.

1428 **G.2 Action Sequencing: Trajectory Error Detection for Missing Step, Additional Step,**  
 1429 **Wrong Temporal Order, Affordance Error**

1430 Action sequencing serves as an intuitive and pragmatic approach to evaluate the effectiveness of LLMs in the  
 1431 context of embodied agents. This task requires LLMs to generate a sequence of executable actions aimed at  
 1432 achieving predefined goals. The evaluation protocol for action sequencing focuses on assessing the LLMs' ability  
 1433 to produce accurate and executable action sequences within embodied environments. The process involves  
 1434 several key steps to ensure that the generated plans are both realistic and effective in achieving the specified  
 1435 goals.

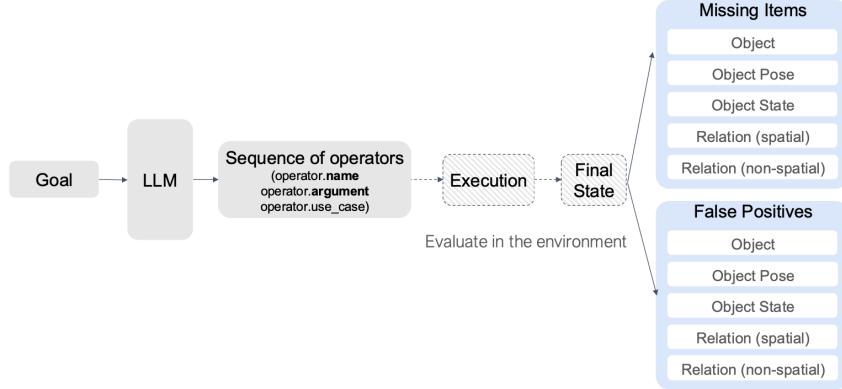


Figure 17: Evaluation pipeline of action sequencing.

1436 **Input Instruction.** In action sequencing tasks, LLMs are given an initial world state  $s_0$ , a goal  $g$ , and  
 1437 general environment information, such as the vocabulary of actions and predicates. The models are required to  
 1438 output a sequence of executable actions  $\bar{a} = \{a_1, a_2, \dots, a_n\}$ , where each action  $a_i$  includes an action name  
 1439 and an ordered list of object names as parameters.

1440 The input includes the objects present in the scene, relative initial state (we use the changes between the initial  
 1441 state and final state to decide a subset of relative objects and their states as the relative initial state), node goals,  
 1442 edge goals, and action goals.

1443 The next step involves rephrasing the concrete goals into natural language, which helps improve the LLM's  
 1444 understanding and execution of the tasks. The rephrased goals are structured as node goals (specifying the state  
 1445 of an object in natural language, such as "*the television is switched on*"), edge goals (describing the relationship  
 1446 between objects in natural language, such as "*the agent is next to the television*"), action goals (describing the  
 1447 required actions in some tasks in natural language, such as "*the robot touches the remote control*").

1448 The evaluation framework for this task is divided into two primary components: trajectory evaluation and goal  
 1449 achievement.

1450 **Trajectory Evaluation.** The trajectory evaluation assesses whether the action sequence  $\bar{a}$  is executable. If  
 1451  $\bar{a}$  is found to be non-executable, errors are categorized into three main classes, each with a fine-grained set of  
 1452 subclasses:

1453 a. **Grammar Errors:**

- 1454 • **Parsing Error:** Evaluates whether the output strictly adheres to specified format requirements.
- 1455 • **Hallucination Error:** The format is correct, but the action names, object names, or predicate  
 1456 names are not in the environment vocabulary.
  - 1457 – *Action Name Hallucination* - Checks for the accuracy of action names supported by the  
 environment.
  - 1459 – *Object Name Hallucination* - Ensures the correctness of object names supported by the  
 environment.
- 1461 • **Action-Argument Number Error** - Verifies that the action or argument parameters are incorrect,  
 1462 mainly if the length of parameters does not meet the specified requirements of the action.

1463 b. **Runtime Errors:**

- 1464 • **Affordance Error:** Determines if the properties of objects allow for the execution of the action.  
 1465 For example, *open(shelf)* is wrong as the shelf cannot be opened.

- 1466     • **Additional Step:** Detects when an action is redundant given the current state. If objects are  
 1467       affordable for the action but the intended effect is already present in the current state, it indicates  
 1468       an unnecessary additional step. For example, *toggle\_on(light)* cannot be executed when the *light*  
 1469       is already *toggled\_on*.  
 1470     • **Missing Step:** Identifies when a required precondition for an action is not met. If (1) properties  
 1471       match, (2) the effect has not been satisfied in the current state, then the execution error stems  
 1472       from an unsatisfied precondition. We then (3) check whether the precondition has never been  
 1473       satisfied in the historical states, which means that a step is missing to trigger the precondition.  
 1474       For example, *release(book)* cannot be executed when the precondition *grasped(book)* is NOT  
 1475       satisfied. If the book has never been grasped by the agent in the historical states, a necessary step  
 1476       *grasp(book)* is missing.  
 1477     • **Wrong Order:** Determines if actions are executed out of sequence. If the precondition is wrong  
 1478       but there has been a precondition being satisfied in the historical states, then the error indicates a  
 1479       wrong order (i.e., the current step should be executed immediately after the relevant historical  
 1480       state). In the previous example, if *grasped(book)* has been previously satisfied (e.g., the book  
 1481       was picked up by the agent but then placed on the table), the order of actions is wrong, and  
 1482       *release(book)* should be promoted to the earlier steps when the book was still grasped.

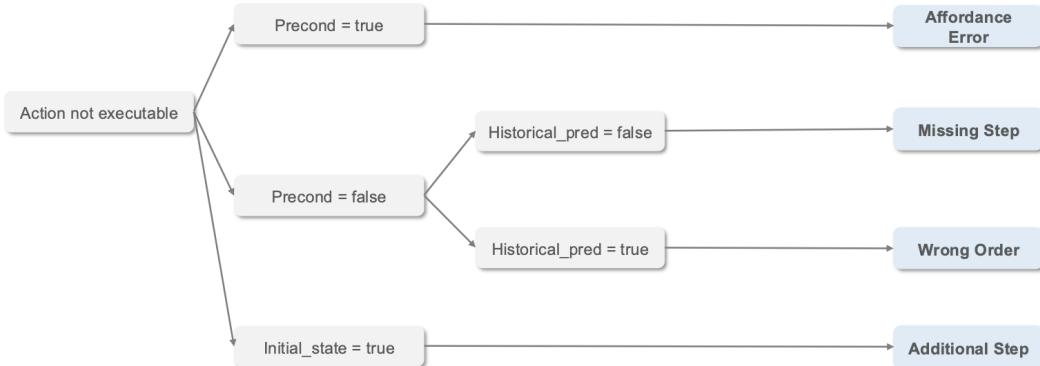


Figure 18: Automatic Error Categorization for Trajectory: (1) Check if the action is affordable based on the properties of objects in the current state. If not, return "Affordance Error". (2) If the action is affordable, check if the intended effect of the action is already satisfied in the current state. If it is, return "Additional Step". (3) If the effect is not redundant, check whether the precondition of the action is satisfied in the current state. If not, proceed to the next step. Otherwise, return "No Runtime Error". (4) If the precondition is not satisfied, check if it has ever been satisfied in any of the historical states. If the precondition has never been satisfied, return "Missing Step". (5) If the precondition has been satisfied in a historical state, return "Wrong Order".

**Trajectory Error Detection:  
Missing Step, Additional Step, Wrong Temporal Order, Affordance Error**

```

def check_runtime_errors(action, current_state, historical_states):
    # Check affordance error
    if not is_affordable(action, current_state):
        return "Affordance Error"

    # Check if action effect is redundant (Additional Step)
    if is_effect_redundant(action, current_state):
        return "Additional Step"

    # Check if precondition for action is met
    if not is_precondition_satisfied(action, current_state):
        # Check if the effect has not been satisfied
        if not is_effect_ever_satisfied(action, historical_states):
            # Check historical states for precondition satisfaction
            if not precondition_satisfied_in_history(action, historical_states):
                return "Missing Step"
            else:
                return "Wrong Order"
        return "No Runtime Error"
  
```

1483

1484 **Goal Satisfaction** The success criteria is based on the accurate achievement of all specified goal conditions  
 1485 and the trajectory to achieve the goal. Both node goals, edge goals, and action goals are checked to ensure that  
 1486 the final state of the environment matches the desired outcomes.

1487 If the action sequence  $\bar{a}$  is executable, the subsequent evaluation examines whether executing  $\bar{a}$  from  $s_0$  satisfies  
 1488 the goal  $g$ . The simulator executes the action sequence  $\bar{a}$  and collects the execution information and intermediate  
 1489 world states  $s_t$  for each action  $a_i \in \bar{a}$ . The execution continues until either an action is not executable (marked  
 1490 as  $a_t$ ) or all actions are executed successfully. The final world state is denoted as  $s_t$ .

1491 Given a goal  $g$ , each goal condition is assigned a category (node goal, edge goal, or action goal) based on a  
 1492 simulator-based classification, similar to Appendix G.1. The satisfaction of each goal condition  $g_i$  by  $s_t$  is then  
 1493 checked using the simulator. The primary metric for this evaluation is the *Success Rate*, calculated as the ratio of  
 1494 tasks where  $g$  is satisfied to the total number of tasks. Additionally, recall is calculated for all goals and different  
 1495 categories of goals to evaluate goal satisfaction.

### 1496 **G.3 Subgoal Decomposition: Converting Subgoal Trajectory to Action Trajectory with BFS 1497 Searching**

1498 The subgoal decomposition module generates a sequence of declarative states. To validate the feasibility of the  
 1499 state transitions, we need to transform the state transitions into an actionable trajectory that can be executed  
 1500 and evaluated in the simulator. As depicted in Figure 15, we address this challenge by utilizing a customized  
 1501 planner to refine the subgoal decomposition output into an actionable sequence. As detailed in Section 2 in the  
 1502 main paper, this subgoal-action mapping function, denoted as  $\mathcal{AM}(\bar{\phi}, s_0)$ , takes the LTL representation of the  
 1503 subgoal sequence  $\bar{\phi}$  and the initial state  $s_0$  as inputs and generates a corresponding state-action sequence  $\bar{t}$ .

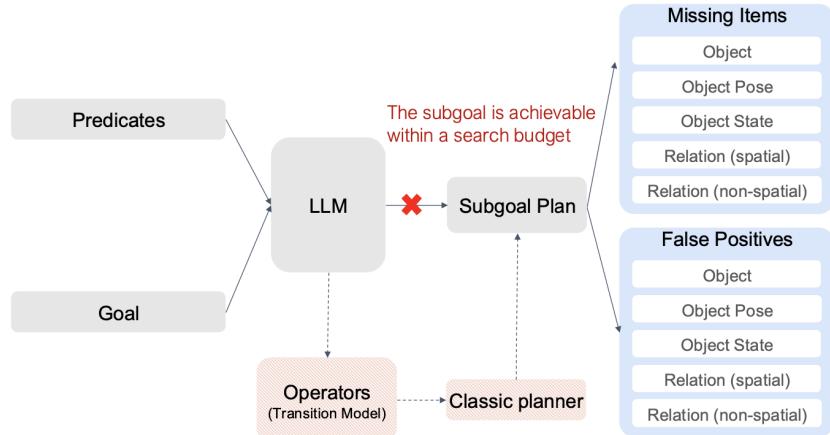


Figure 19: Evaluation pipeline of subgoal decomposition. The subgoal decomposition module generates a sequence of declarative states, which need to be transformed into an actionable trajectory for execution in the simulator.

1504 We implement this mapping function using a Breadth-First Search (BFS) algorithm. We show this process in  
 1505 Figure 19, where we find a sequence of actions that can transition the agent from the initial state to the desired  
 1506 subgoal states.

1507 The BFS algorithm starts from the initial state  $s_0$  and expands the search frontier by exploring all possible  
 1508 actions at each step. It maintains a queue of states to be visited and keeps track of the path from the initial state  
 1509 to each visited state. The search continues until a state satisfying the first subgoal is reached. The process is then  
 1510 repeated, using the reached subgoal state as the new initial state and the next subgoal as the target.

1511 The subgoal-action mapping function  $\mathcal{AM}(\bar{\phi}, s_0)$  takes the LTL representation of the subgoal sequence  $\bar{\phi}$  and  
 1512 the initial state  $s_0$  as inputs and returns a state-action sequence  $\bar{t}$  that achieves the subgoals. The function can be  
 1513 described as follows:

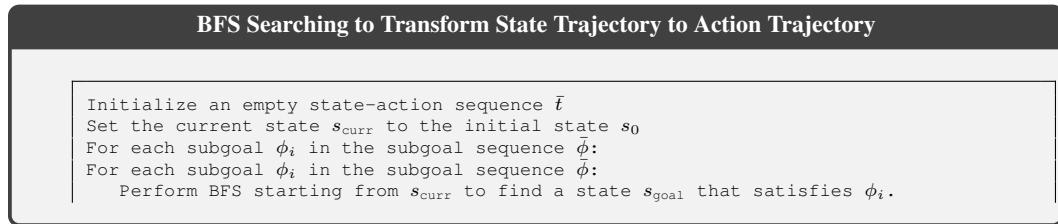




Figure 20: Transition Model evaluation metrics based on Bipartite Graph Matching for pre-conditions and post-effects.

1515

Extract the path from  $s_{curr}$  to  $s_{goal}$  and append the corresponding state-action pairs to  $\bar{t}$ .  
Set  $s_{curr}$  to  $s_{goal}$ .  
Return the complete state-action sequence  $\bar{t}$ .

1516

By using the BFS algorithm to find action sequences that connect the subgoal states, we can effectively transform the declarative subgoal decomposition into an imperative action trajectory. This allows us to evaluate the subgoal decomposition module in the simulator and assess its effectiveness in guiding the agent towards the desired goal.

1519

Note that the BFS algorithm guarantees finding the shortest path between subgoals if one exists, but it may be computationally expensive in large state spaces. In our current benchmark, due to limited pre-defined action space and state space (detailed in Appendix Q.2 and Appendix Q.1), we do not need to control the searching space. In practice, with larger action space and state space, heuristics or domain-specific knowledge can be incorporated to guide the search and improve efficiency.

1524

By employing this approach, we can effectively evaluate the subgoal decomposition module within the simulator, ensuring that the generated subgoals can be successfully grounded and executed. This allows us to assess the quality of the subgoal decomposition and its impact on the overall decision-making process, even though the direct output of the module is declarative states rather than imperative actions.

1528

#### G.4 Transition Modeling: Evaluating with PDDL Planners

1529

The evaluation protocol for transition modeling aims to rigorously assess the ability of Large Language Models (LLMs) to predict accurate action preconditions and effects within embodied agent scenarios. This evaluation is critical for understanding the proficiency of LLMs in modeling the dynamics of physical interactions in simulated environments.

1533

The process begins with the preparation of input data, which includes goals, scripts, initial states, and final states provided in the dataset. Initially, predicates representing relations and properties within the domain are annotated. Based on these annotated predicates, PDDL (Planning Domain Definition Language) operators are defined. Subsequently, a comprehensive PDDL domain file is constructed, encapsulating the full set of predicates, including states, relations, and properties necessary for the planning tasks. Next, relevant objects are identified from the script, initial state, and final states. These objects are used to ground the goals into specific goals for each task. The problem file is formulated using the relevant objects, initial states, and grounded goals.

1540 With the domain and problem files prepared, the LLM is tasked with predicting the transition model, that is, the  
1541 action preconditions and effects. The inputs to the LLM include the predicate list, problem file, action name, and  
1542 action parameters, and the output is the actions body predicted by LLM.

1543 The output evaluation involves several key steps to assess the accuracy and reliability of the predicted transition  
1544 model. First, the predicted preconditions and effects are extracted from the LLM-generated actions. These  
1545 predicates are categorized to facilitate detailed analysis of the model's performance. The evaluation protocol  
1546 uses several metrics to assess the performance of the LLMs. The logical scoring function evaluates the similarity  
1547 between the predicted and gold action sequences by comparing the logical structure of preconditions and effects.  
1548 The success rate by planner measures the proportion of correctly predicted actions that exists a feasible solution  
1549 to achieve the desired goals under PDDL planner[70]<sup>‡</sup>. Sensitivity analysis demonstrates how difficult it is for  
1550 LLM to predict specific actions in different tasks.

### Logic Form Accuracy function(expr1, expr2)

```
if expr1 and expr2 are literals:  
    return expr1 == expr2  
if type(expr1) != type(expr2):  
    return 0  
if isinstance(expr1, Not):  
    return match_expressions(expr1.expression, expr2.expression)  
if isinstance(expr1, When):  
    if match_expressions(expr1.condition, expr2.condition)  
        and match_expressions(expr1.consequence, expr2.consequence):  
            return 1  
    return 0  
if isinstance(expr1, (Exists, Forall)):  
    if quantifier are the same:  
        return match_expressions(expr1.body, expr2.body)  
    return 0  
if isinstance(expr1, (And, Or)):  
    adj_matrix = np.zeros((len(expr1.args), len(expr2.args)), dtype=int)  
    for i in range(size1):  
        for j in range(size2):  
            adj_matrix[i, j] = match_expressions(sub1[i], sub2[j])  
    match_result = maximum_bipartite_matching(sparse_matrix, perm_type='column')  
    total_match = sum(adj_matrix[i, match_result[i]]  
    for i in range(size1) if match_result[i] != -1)  
    max_possible_match = min(size1, size2)  
    return total_match / max_possible_match if max_possible_match > 0 else 0
```

1551

1552 **Logic Form Accuracy.** In PDDL, logical connectives such as *and*, *or*, *not*, *when*, *forall*, and *exists* are used to  
1553 define preconditions and effects. The logical matching evaluates the similarity between the predicted and ground  
1554 truth preconditions or effects by parsing them into clauses of disjunctive logical form and performing bipartite  
1555 matching.

1556 Given a set of predicted clauses  $P = \{p_1, p_2, \dots, p_n\}$  and a set of ground truth clauses  $G = \{g_1, g_2, \dots, g_m\}$ ,  
1557 we define an adjacency matrix  $A$  where  $A[i, j] = \text{match}(p_i, g_j)$  represents the similarity between the  $i$ -th  
1558 predicted clause and the  $j$ -th ground truth clause. The function  $\text{match}(p, g)$  returns 1 if the clauses are identical,  
1559 and 0 otherwise. If the ground truth clause  $g_j$  is empty,  $\text{match}(p, g) = 1$  if and only if  $p_i$  is empty as well.

1560 For clauses containing connectives, the clauses are recursively expanded and bipartite matching is conducted  
1561 on the expanded nodes. If  $p$  and  $g$  are literals, the match function compares their equivalence directly. If  
1562 either clause contains a connective, the clause is further decomposed and bipartite matching is applied to the  
1563 sub-clauses.

1564 The evaluation metrics for logical matching include precision, recall, and F1-score. True Positives (TP) are the  
1565 matched clauses, False Positives (FP) are the unmatched clauses in the predicted set, and False Negatives (FN)  
1566 are the unmatched clauses in the ground truth set.

1567 The process begins by parsing the preconditions or effects into disjunctive logical form and constructing the  
1568 adjacency matrix  $A$  based on clause similarity. Bipartite matching is performed to find the optimal pairing of  
1569 predicted and ground truth clauses. The matching process is recursive, expanding nested clauses and applying  
1570 bipartite matching at each level. If a sub-clause does not match, the entire clause is considered a non-match.

1571 Finally, the precision, recall, and F1-score are calculated using the formulas above, providing a robust measure  
1572 of the LLM's accuracy in modeling logical relationships within PDDL tasks.

<sup>‡</sup>[https://github.com/ronuchit/pddlgym\\_planners](https://github.com/ronuchit/pddlgym_planners)

1573 **Success Rate by External Planners.** The planner success rate evaluates the success of the LLM in generating  
1574 correct and executable action sequences for each task category. The overall success rate across all categories is  
1575 also reported to provide a comprehensive assessment.

1576 Given a set of programs  $P = \{p_1, p_2, \dots, p_n\}$  and a corresponding set of categories  $C = \{c_1, c_2, \dots, c_k\}$ , let  
1577  $P_{c_j} \subseteq P$  denote the subset of programs belonging to category  $c_j$ . The success rate for category  $c_j$  is defined as:

$$\text{SuccessRate}(c_j) = \frac{\sum_{p \in P_{c_j}} \text{IsSuccessful}(p)}{|P_{c_j}|}$$

1578 where  $\text{IsSuccessful}(p)$  is a binary function that returns 1 if the planner successfully generates an executable  
1579 action sequence for program  $p$ , and 0 otherwise.

1580 The overall success rate across all categories is given by:

$$\text{OverallSuccessRate} = \frac{\sum_{j=1}^k \sum_{p \in P_{c_j}} \text{IsSuccessful}(p)}{\sum_{j=1}^k |P_{c_j}|}$$

## 1581 H Full Results with 15 models

1582 In this section, we provide the full results and analysis for each ability module using 15 different models, and  
1583 provide a further analysis on the potential coorelation with factors in Appendix H.5.

### 1584 H.1 Goal Interpretation

1585 Table 10 examines the performance of various LLMs in translating natural language task instructions into  
1586 actionable symbolic goals within two different simulators: VirtualHome (V) and BEHAVIOR (B). Additionally,  
1587 Figures 21 and 22 show detailed qualitative error cases based on error type and goal type. Below are the detailed  
1588 result analysis and error case studies.

#### 1589 H.1.1 Result Analysis

1590 Table 10 contains the quantitative benchmarking results of the Goal Interpretation ability module for both  
1591 Behavior and VirtualHome simulators. The results are mainly broken down into three categories: State Goals or  
1592 Unary Goals (which describe the goal state of just one object), Spatial Goals or Binary Goals (which describe  
1593 the spatial relationship goal of two interacting objects), and Action Goals specific to the VirtualHome simulator  
1594 (which describe the actions that an agent should finish to complete the task). Finally, we include an Overall  
1595 section, where all goals' performances are combined with a micro average.

1596 In Table 10, our evaluation metrics consist of precision, recall, and F1 score for each goal type and overall  
1597 performance. In the case of the goal interpretation task, precision measures how unlikely a model is to make  
1598 false positive goal predictions given reasonable instructions and task information, whereas recall measures how  
1599 unlikely a model is to leave out ground truth goals in their prediction. The F1 score combines these two metrics  
1600 and reflects overall model performance for each category.

1601 Based on the results in Table 10, Gemini 1.5 Pro achieves top overall goal interpretation performance (F1 score)  
1602 in both VirtualHome and BEHAVIOR simulators over other LLMs, whereas Claude-3 Opus achieves the highest  
1603 successful ground truth goal retrieval rate (Recall) in both simulators. GPT-4o also demonstrates very high  
1604 performance in both simulators with very close precision, recall, and F1 scores compared to the top performer on  
1605 each metric. Among open-source models, the top performer by a significant margin is the Llama-3 70B Instruct  
1606 model, leading in overall F1 scores for both simulators.

1607 For individual goal types, the performance trends are less uniform: GPT-4o demonstrates a very strong spatial  
1608 reasoning ability, reaching top overall F1 score in both Behavior and VirtualHome's spatial goal categories.  
1609 Gemini 1.5 Pro leads in terms of overall performance in VirtualHome action goals and Behavior state goals,  
1610 while the Cohere Command R model leads in terms of VirtualHome state goals.

#### 1611 H.1.2 Case Studies

1612 In Figure 21, we show detailed examples of Behavior and VirtualHome goal interpretation errors broken down  
1613 into two main categories: goal prediction errors and grammar errors.

1614 Grammar errors can be further broken down into hallucination errors (such as object/state/action hallucination)  
1615 and format errors (such as JSON parsing errors). In our empirical experiments, we find that given a reasonably  
1616 clear prompt with in-context examples, commercial LLMs such as the GPT and Claude family models are

1617 extremely unlikely to display format errors, while open-source models such as Llama3 8B and Mixtral 8x22B  
 1618 tend to make a small number of mistakes.

1619 Goal prediction errors can be further broken down into missing goals and false positive goals, which can be  
 1620 quantitatively measured by recall and precision scores, respectively. Empirically, we find two kinds of errors to  
 1621 be common among both open-source and proprietary LLMs:

- 1622 (1) As shown in the false positive goal error sections of Figure 21 (a) and (b), LLMs tend to output  
 1623 intermediate goal states in place of final goal predictions. This phenomenon can be intuitively  
 1624 explained by the LLMs’ reliance on Chain-of-Thought Reasoning [71] for question answering. While  
 1625 our carefully designed system and user prompts restrict LLMs to directly output the predicted goals  
 1626 without any explicit intermediate explanation, models may still try to perform such intermediate  
 1627 reasoning steps in their structured output, thus leading to such errors.
- 1628 (2) As shown in the missing goal error sections of Figure 21 (a) and (b), LLMs tend to leave out simple  
 1629 object spatial relationships in their output. For the task “serving a meal,” with ground truth goal  
 1630 condition  $\text{ONTOP}(\text{chicken.0}, \text{plate.2})$  and  $\text{ONTOP}(\text{plate.2}, \text{table.1})$ , GPT-4o (along with many other  
 1631 models) mistakenly predicts  $\text{ONTOP}(\text{chicken.0}, \text{table.1})$ . While the expression “chicken ontop of table”  
 1632 is acceptable in a conversational setting, it presents a completely wrong set of physical relationships  
 1633 between the objects chicken, plate, and table. This type of error is common in both Behavior and  
 1634 VirtualHome simulators for various models, highlighting a significant issue of imprecise spatial  
 1635 relationship description in applying LLMs for embodied planning and robotic control, where physical  
 1636 precision is crucial for task success.

Table 10: All goal evaluation results (%) for goal interpretation

Model Name	Goal Interpretation																							
	State			Spatial			Action			Overall														
	Precision		Recall	Precision		Recall	Precision		Recall	Precision		Recall	Precision		Recall	F1								
V	B	V	B	V	B	V	B	V	B	V	B	V	B	V	B	V	B							
Claude-3 Haiku	21.8	22.8	58.9	93.5	31.8	36.7	24.2	64.5	50.8	64.6	32.8	64.6	12.2	-	<b>95.7</b>	-	21.6	-	18.0	41.5	63.2	71.2	28.0	52.5
Claude-3 Sonnet	23.3	36.8	57.1	88.9	33.1	52.0	26.6	76.2	<b>53.0</b>	<b>79.8</b>	35.5	77.9	12.4	-	85.8	-	21.7	-	19.3	60.2	61.5	81.9	29.4	69.4
Claude-3 Opus	27.0	72.6	66.9	93.5	38.5	81.7	22.6	75.2	46.8	79.2	30.5	77.1	14.5	-	92.6	-	25.1	-	20.7	72.2	<b>65.0</b>	<b>82.5</b>	31.4	77.0
Cohere Command R	<b>51.1</b>	7.7	<b>69.6</b>	31.4	<b>58.9</b>	12.4	34.5	56.8	21.3	55.0	26.3	55.9	3.6	-	38.9	-	6.5	-	27.4	28.2	55.7	49.6	36.7	36.0
Cohere Command R+	20.9	23.3	52.0	79.1	29.8	36.0	17.9	66.7	15.2	61.5	16.4	64.0	10.4	-	82.6	-	18.5	-	14.9	42.0	44.5	65.5	22.4	51.2
Gemini 1.0 Pro	25.3	27.4	57.9	81.1	34.9	41.0	17.0	75.2	20.6	70.4	18.6	72.7	9.9	-	68.7	-	17.2	-	16.2	51.0	45.2	72.8	23.8	60.0
Gemini 1.5 Flash	23.6	55.8	57.9	94.1	33.5	70.1	19.8	76.6	21.1	76.7	20.5	76.7	13.5	-	90.1	-	23.5	-	18.2	69.7	50.8	80.7	26.8	74.8
Gemini 1.5 Pro	45.4	<b>94.0</b>	49.1	92.8	47.2	<b>93.4</b>	<b>40.0</b>	74.4	9.7	76.7	15.6	75.6	<b>26.8</b>	-	80.9	-	<b>40.3</b>	-	<b>35.2</b>	<b>78.8</b>	41.1	80.4	<b>37.9</b>	<b>79.6</b>
GPT-3.5-turbo	22.4	52.0	50.0	66.7	30.9	58.5	8.5	51.5	18.8	46.9	11.7	49.1	15.2	-	60.5	-	24.4	-	15.7	49.5	40.5	51.4	22.7	50.4
GPT-4-turbo	28.6	70.4	58.5	86.9	38.4	77.8	24.7	77.5	32.9	76.4	28.2	76.9	19.0	-	82.1	-	30.9	-	24.0	75.6	53.8	78.8	33.2	77.2
GPT-4o	29.0	67.1	60.0	94.8	39.1	78.6	31.5	<b>81.1</b>	43.6	78.5	<b>36.6</b>	<b>79.8</b>	20.5	-	85.8	-	33.1	-	26.4	76.5	59.1	82.2	36.5	79.2
Llama3 8B Instruct	21.7	17.3	54.4	80.4	31.0	28.4	14.0	51.4	7.4	20.8	9.7	29.6	11.1	-	79.4	-	19.4	-	15.5	24.1	41.9	34.3	22.6	28.3
Llama3 70B Instruct	23.9	69.5	61.2	<b>95.4</b>	34.3	80.4	22.6	70.0	37.5	73.3	28.2	71.6	11.2	-	88.8	-	19.8	-	17.5	64.7	58.0	78.3	26.9	70.9
Mistral Large	23.6	63.5	59.1	92.2	32.8	75.2	23.7	75.1	40.3	76.2	29.8	75.6	11.2	-	84.0	-	19.7	-	17.5	69.6	57.1	79.8	26.8	74.3
Mixtral 8x22B MoE	23.6	22.9	56.9	83.7	33.4	36.0	22.2	70.7	36.3	67.7	27.5	69.2	11.2	-	<b>94.8</b>	-	20.0	-	17.4	44.4	56.2	71.3	26.6	54.7

### Takeaways of Goal Interpretation

- (1) **Top-performing Models:** Gemini 1.5 Pro achieves top overall goal interpretation performance (F1-score) in both VirtualHome and BEHAVIOR simulators over other LLMs. Whereas Claude-3 Opus achieves the highest successful ground truth goal retrieval rate (Recall) in both simulators.
- (2) **Grammar Errors:** Current SOTA proprietary LLMs generally make very little to no grammar errors in goal interpretation, whereas top open-source LLMs including Llama3 70B Instruct tend to suffer more from format/parsing errors and object/state hallucination.
- (3) **Goal Prediction Errors:** Both open-source and proprietary LLMs suffer significantly from two common types of goal prediction errors:
  - (a) Generating intermediate goals in place of final goals (when asked to directly output final goals), resulting in false positive goal predictions.
  - (b) Omitting conversationally uncommon spatial relationship goals that are essential for precise robotic control, resulting in missing goals (reporting bias).

1637

## H.2 Subgoal Decomposition

### H.2.1 Result Analysis

1640 Table 11 examines the performance of various LLMs in decomposing high-level goals into actionable subgoals  
 1641 within two different simulators: VirtualHome (V) and BEHAVIOR (B). Additionally, Table 12 investigates  
 1642 fine-grained goal satisfaction for each LLM. Below are the results and related error analysis.

VirtualHome: Goal Interpretation – Errors			
Goal Prediction Errors		Grammar Errors	
<b>Missing Goal</b> Model: Claude-3 Haiku Task Name: Cook Some Food	<b>False Positive</b> Model: GPT-4o Task Name: Drink	<b>Hallucination</b> Model: Llama3 70B Instruct Task Name: Drink	<b>Format Error</b> Model: Mixtral 8x22b MoE Task Name: Work
<b>Ground Truth Goal</b> <code>on(open.0)</code> <code>closed(open.0)</code>	<b>Ground Truth Goal</b> <code>empty(glass.0)</code> <code>open(cupboard.0)</code>	<b>Ground Truth Goal</b> <code>or holds_right_hand (cup.1)</code> <code>holds_left_hand (cup.1)</code>	<b>Ground Truth Goal</b> <code>node_goals: [sitting(character.0)...</code>
<b>LLM Output</b> ... <code>plugged_in(open.0)</code> <code>on(open.0)</code>	<b>LLM Output</b> <code>full(glass.0)</code> <code>open(cupboard.0)</code>	<b>LLM Output</b> ... <code>inside (cup.1, hands_both.0)</code>	<b>LLM Output</b> ... <code>SYMBOLIC NODE GOALS:\n[...</code>
<b>Error Reason</b> ✖ Reporting Bias, Leaving Out Goal: <code>closed(open.0)</code>	<b>Error Reason</b> ✖ Outputting Intermediate Goal: <code>full(glass.0)</code>	<b>Error Reason</b> ✖ Object Hallucination: “hands_both.0” is not a valid VirtualHome simulator object	<b>Error Reason</b> ✖ Wrong JSON Key: SYMBOLIC NODE GOALS

(a) VirtualHome

BEHAVIOR: Goal Interpretation – Errors			
Goal Prediction Errors		Grammar Errors	
<b>Missing Goal</b> Model: GPT-4o Task Name: bringing_in_wood	<b>False Positive</b> Model: Claude-3 Sonnet Task Name: bottling_fruit	<b>Hallucination</b> Model: GPT-4-turbo Task Name: cleaning_up_after_meal	<b>Format Error</b> Model: GPT-3.5-turbo Task Name: bottling_fruit
<b>Ground Truth Goal</b> <code>onfloor(plywood.4,floor.2)</code> <code>onfloor(plywood.3,floor.1)</code>	<b>Ground Truth Goal</b> <code>closed(jar.1)</code> <code>inside(peach.0,jar.1)</code>	<b>Ground Truth Goal</b> <code>not stained(tray.2)</code> <code>not stained(bowl.1)</code>	<b>Ground Truth Goal</b> <code>inside(peach.0,jar.1)</code> <code>closed(jar.1)</code>
<b>LLM Output</b> ... <code>inside(plywood.4, room.1)</code> <code>inside(plywood.3, room.1)</code>	<b>LLM Output</b> ... <code>open(jar.1)</code> <code>inside(peach.0, jar.1)</code>	<b>LLM Output</b> ... <code>cleaned(tray.2)</code> <code>bowl(bowl.1)</code>	<b>LLM Output</b> ... <code>inside(peach.0, jar.1)</code> <code>closed(jar.1)</code>
<b>Error Reason</b> ✖ Reporting Bias, Leaving Out Goals: <code>onfloor(plywood.4,floor.2)</code> <code>onfloor(plywood.3,floor.1)</code>	<b>Error Reason</b> ✖ Outputting Intermediate Goal: <code>open(jar.1)</code>	<b>Error Reason</b> ✖ State Hallucination: “cleaned” is not a valid state of the Behavior simulator	<b>Error Reason</b> ✖ Additional Parentheses: <code>closed(jar.1))</code>

(b) BEHAVIOR

Figure 21: Goal evaluation error examples for goal interpretation.

VirtualHome: Goal Interpretation – Errors		BEHAVIOR: Goal Interpretation – Errors		
State Goal Error	Relation Goal Error	Action Goal Error	State Goal Error	Relation Goal Error
Model: GPT-4o Task Name: Turn on Light	Model: Claude-3 Opus Task Name: Wash Teeth	Model: Claude-3 Haiku Task Name: Change TV Channel	Model: Mistral Large Task Name: storing the groceries	Model: GPT-4o Task Name: serving_a_meal
<b>Ground Truth Goal</b> <code>on(light.0)</code> <code>plugged_in(light.0)</code>	<b>Ground Truth Goal</b> <code>hold_right_hand(tooth_brush.0)</code>	<b>Ground Truth Goal</b> <code>touch(remote_control.0)</code> <code>push(remote_control.0)</code>	<b>Ground Truth Goal</b> <code>closed(fridge.0)</code> <code>inside(lettuce.1, fridge.0)</code>	<b>Ground Truth Goal</b> <code>ontop(chicken.0, plate.2)</code> <code>ontop(plate.2, table.1)</code>
<b>LLM Output</b> ... <code>on(light.0)</code>	<b>LLM Output</b> ... <code>hold_right_hand(tooth_paste.0)</code>	<b>LLM Output</b> <code>grab(remote_control.0)</code> <code>switchon(remote_control.0)</code>	<b>LLM Output</b> <code>inside(lettuce.1, fridge.0)</code> <code>open(fridge.0)</code>	<b>LLM Output</b> ... <code>ontop(chicken.0, table.1)</code>
<b>Error Reason</b> ✖ Missing State Goal: <code>plugged_in(light.0)</code>	<b>Error Reason</b> ✖ Wrong Relation Goal: <code>hold_right_hand(tooth_paste.0)</code>	<b>Error Reason</b> ✖ Wrong Action Goals: <code>grab</code> and <code>switchon</code>	<b>Error Reason</b> ✖ Wrong State Goal: <code>open(fridge.0)</code>	<b>Error Reason</b> ✖ Wrong Relation Goal: <code>ontop(chicken.0, table.1)</code>

(a) VirtualHome

(b) BEHAVIOR

Figure 22: Goal satisfaction error examples for goal interpretation

1643 **VirtualHome (1) Executable and Goal Success Rate.** The top-performing models are Gemini 1.5 Flash  
1644 and GPT-4o, with success rates of 89.05% and 88.76%, respectively. For executable side, GPT-4-turbo and  
1645 Gemini 1.5 Flash achieves the best executable performance, with the rates of 94.08%. Surprisingly, despite  
1646 designed smaller than Gemini 1.5 Pro, Gemini 1.5 Flash achieves the best performance overall, suggesting robust  
1647 subgoal decomposition that aligns well with the VirtualHome simulator’s requirements. Conversely, models like  
1648 Llama3 8B, with a success rate of 48.82%, illustrate the challenges of effective subgoal decomposition.

1649 **(2) Grammar Errors.** The occurrence of grammar errors is notably low across the board. Specifically, most  
1650 of SOTA LLMs make no errors in parsing and predicate arguments number. This indicates that current SOTA  
1651 LLMs can follow the subgoal syntax. However, some models are prone to hallucinate non-existing predicates.  
1652 Specifically, GPT-4o tends to hallucinate the action *POUR* when dealing with the task ‘make coffee’, which is  
1653 not defined in the subgoal decomposition setting.

1654 **(3) Runtime Errors.** In terms of runtime errors, compared with wrong temporal order and affordance errors,  
1655 LLMs are prone to make missing steps and additional steps errors, while additional steps errors are the  
1656 most critical. (1) Wrong temporal order. Closed-sourced LLMs perform well in understanding the temporal  
1657 requirement among subgoals, while open-sourced LLMs like Llama3 8B will make more temporal order errors.  
1658 For the wrong cases, LLMs tend to ignore the sitting or lying state of the agent and fail to call action *STANDUP*  
1659 before they applying some actions requiring standing state, yet this state has been achieved earlier. (2) Missing  
1660 step. LLMs sometimes are prone to fail satisfy some preconditions when applying actions. Among all LLMs,  
1661 GPT-3.5-turbo performs worst in this error type. Specifically, it tends to ignore opening a closed object before  
1662 fetching something inside it. (3) Affordance. Overall, Llama3 8B performs much worse than other LLMs,  
1663 meaning it cannot understand the semantics very well. (4) Additional steps. All LLMs are prone to produce  
1664 additional steps errors, even for SOTA LLMs like GPT-4o and Claude-3 Opus. This is mainly because in the  
1665 initial scene state, some of the goals have already been achieved, yet LLMs still prefer to plan the satisfied goals  
1666 in its output.

1667 **(4) Goal Satisfaction Analysis.** In Table 12, we can see that LLMs perform well in understanding and satisfying  
1668 state goals. This is because states are usually isolated from other objects with less logical requirement. For  
1669 example, unary states like *on* and *off* can be easily achieved by calling *SWITCH ON* and *SWITCH OFF*,  
1670 respectively. Also, unary states are less like to be affected by other objects states changing in VirtualHome  
1671 . However, a relation usually involves two objects, which increases the difficulty to satisfy the goals. For  
1672 example, relation like *ontop(agent, chair)* requires that the agent should be standing and the chair should be  
1673 empty. Furthermore, LLMs performance varies across action goal success rate. GPT-4o achieves an accuracy of  
1674 93.21%, while Cohere Command R and Llama3 8B achieves less than 70% accuracy. An possible explanation  
1675 is that stronger LLMs can better understand the semantics of actions, while weaker LLMs have trouble in  
1676 understanding the difference between actions and states. Generally, most LLMs are capable to achieve more than  
1677 80% goals defined in our annotated data. This could be because that goals in RobotHow are relatively simple  
1678 and straightforward.

1679 **BEHAVIOR (1) Executable and Goal Success Rate.** GPT-4o stands out with the highest success rate  
1680 at 48.00% and an executable rate 55.00%, followed by Claude-3 Opus with a success rate of 34.00% and an  
1681 executable rate of 42.00%. Overall, all LLMs fail to achieve a high performance in BEHAVIOR . This could be  
1682 the tasks representation of BEHAVIOR is much more complicated than that of VirtualHome . For example, most  
1683 tasks in BEHAVIOR have quantifiers like *forall* and *forpairs* with complex spatial or temporal requirement, yet  
1684 VirtualHome tasks have much easier goal definitions.

1685 **(2) Grammar Errors.** Most of LLMs are proficient in generating grammatically correct subgoal plans does has  
1686 no parsing errors and correct number of predicate parameters. There are a few exceptions, though. One-fifth of  
1687 cases are parsing errors for Cohere Command R, and most of them are illegal tokens or syntax for LTL parser  
1688 like ‘[’. Additionally, LLMs like Cohere Command R and Llama3 8B tend to hallucinate non-existing objects in  
1689 the scene.

1690 **(3) Runtime Errors.** Most runtime errors occurring in BEHAVIOR across various LLMs are due to missing  
1691 steps. In fact, over half of the total cases involve such errors in the majority of LLMs. Even GPT-4o, which  
1692 has the lowest error rate for missing steps, encounters these errors in 37% of cases. This can be attributed  
1693 primarily to two reasons. First, the precondition checking in BEHAVIOR is stricter than that in VirtualHome .  
1694 For instance, VirtualHome does not verify preconditions such as whether an agent is holding a cleaning tool  
1695 before executing the *WASH* action. In contrast, BEHAVIOR requires such conditions to be satisfied before  
1696 invoking a similar action like *CLEAN*. Secondly, the complexity of tasks in BEHAVIOR often leads LLMs to  
1697 overlook seemingly trivial actions, such as opening a closed container before retrieving an item inside it. This  
1698 issue persists even when a heads-up is included in the prompt. Interestingly, LLMs tend to make significantly  
1699 fewer additional step errors in BEHAVIOR compared to VirtualHome . This is likely because the majority of  
1700 task goals in BEHAVIOR are not satisfied in the initial state, making additional steps less frequent.

1701 **(4) Goal Satisfaction Analysis.** From the statistics, it is evident that LLMs do not perform well in achieving  
1702 both state goals and relation goals. The discrepancy arises because we use different metrics to evaluate state  
1703 and relation goals as opposed to those used in VirtualHome . In BEHAVIOR , most state and relation goals  
1704 are encapsulated within quantifiers. Consequently, quantifiers such as *forall* or *forpairs* tend to fail if even a  
1705 single state or relation goal is not met. Furthermore, since most quantifiers can have multiple solutions, it is  
1706 challenging to obtain accurate statistics for state and relation goals within quantifiers. Furthermore, our metric  
1707 considers quantifiers as a combined entity for both state and relation goals. This approach also contributes to the  
1708 low success rates for both state and relation goals.

## 1709 **H.2.2 Case Studies**

1710 In this section, we would like to investigate several error examples made in subgoal decomposition by LLMs.  
1711 Specifically, we explore in both VirtualHome and BEHAVIOR .

Table 11: All trajectory evaluation results (%) for subgoal decomposition.

Model	Goal Evaluation				Trajectory Evaluation							
	Goal SR		Execution SR		Grammar Error ( $\downarrow$ )				Runtime Error ( $\downarrow$ )			
	V	B	V	B	Parsing	Hallucination	Action-Arg Num	Wrong Order	Missing Step	Affordance	Additional Step	
Claude-3 Haiku	78.4	29.0	82.8	35.0	0.3	<b>0.0</b>	2.4	1.0	1.8	0.0	1.8	2.0
Claude-3 Sonnet	83.1	38.0	86.4	43.0	<b>0.0</b>	2.0	1.8	<b>0.0</b>	<b>0.0</b>	2.0	0.6	3.0
Claude-3 Opus	87.0	39.0	90.0	47.0	0.3	<b>0.0</b>	3.6	3.0	<b>0.0</b>	<b>0.0</b>	1.2	5.0
Gemini 1.0 Pro	70.4	23.0	84.6	33.0	0.6	2.0	3.3	4.0	2.4	<b>0.0</b>	1.2	3.0
Gemini 1.5 Flash	<b>89.1</b>	34.0	<b>94.1</b>	42.0	<b>0.0</b>	2.0	1.5	1.0	<b>0.0</b>	<b>0.0</b>	0.6	2.0
Gemini 1.5 Pro	87.0	31.0	91.1	37.0	<b>0.0</b>	1.0	1.5	0.0	1.8	1.0	<b>0.0</b>	3.0
GPT-3.5-turbo	69.2	24.0	81.4	36.0	1.5	2.0	<b>0.0</b>	3.0	0.6	<b>0.0</b>	1.5	3.0
GPT-4-turbo	85.5	37.0	<b>94.1</b>	47.0	<b>0.0</b>	<b>0.0</b>	1.8	3.0	<b>0.0</b>	<b>0.0</b>	1.5	9.0
GPT-4o	88.8	<b>48.0</b>	90.2	<b>55.0</b>	<b>0.0</b>	<b>0.0</b>	6.2	3.0	<b>0.0</b>	<b>0.0</b>	1.2	5.0
Cohere Command R	71.3	15.0	79.6	25.0	2.1	22.0	3.9	11.0	0.9	<b>0.0</b>	1.5	<b>0.0</b>
Cohere Command R+	79.0	24.0	83.7	37.0	1.5	2.0	4.5	4.0	2.1	<b>0.0</b>	0.9	5.0
Mistral Large	84.3	30.0	92.0	38.0	0.3	1.0	1.8	3.0	0.3	<b>0.0</b>	2.1	4.0
Mixtral 8x22B MoE	80.5	27.0	90.2	33.0	0.3	0.0	2.4	4.0	<b>0.0</b>	<b>0.0</b>	3.0	2.0
Llama3 8B	48.8	21.0	58.0	29.0	0.6	2.0	2.4	11.0	0.6	<b>0.0</b>	6.8	6.0
Llama3 70B	78.4	20.0	87.3	30.0	<b>0.0</b>	1.0	2.4	5.0	0.9	1.0	2.4	8.0

Table 12: All goal success results (%) for action sequencing and subgoal decomposition.

Model	Action Sequencing						Subgoal Decomposition					
	State Goal		Relation Goal		Action Goal		State Goal		Relation Goal		Action Goal	
	V	B	V	B	V	B	V	B	V	B	V	B
Claude-3 Haiku	59.4	27.0	42.8	38.7	87.8	-	61.4	35.5	89.4	26.0	82.2	34.8
Claude-3 Sonnet	79.5	41.0	67.2	<b>59.8</b>	85.1	-	77.2	<b>54.6</b>	89.1	37.0	<b>89.3</b>	49.8
Claude-3 Opus	63.7	45.0	71.1	53.0	77.0	-	69.1	50.8	92.4	43.0	88.6	41.6
Gemini 1.0 Pro	52.5	28.0	32.8	32.0	77.0	-	52.6	30.9	84.4	26.0	61.5	31.1
Gemini 1.5 Flash	79.5	34.0	63.3	50.0	88.5	-	76.9	45.6	<b>93.5</b>	44.0	88.3	36.0
Gemini 1.5 Pro	81.7	41.0	76.7	43.2	89.2	-	82.0	42.6	91.2	31.0	72.5	37.1
GPT-3.5-turbo	29.1	20.0	15.6	22.6	64.2	-	33.7	21.9	84.7	28.0	54.4	28.5
GPT-4-turbo	74.8	39.0	72.2	39.5	89.9	-	77.7	39.3	<b>93.5</b>	45.0	84.2	46.1
GPT-4o	<b>83.1</b>	<b>49.0</b>	71.1	45.5	89.9	-	81.2	46.5	92.1	<b>50.0</b>	84.2	<b>53.2</b>
Cohere Command R	18.4	20.0	31.1	25.9	48.0	-	29.4	24.3	85.3	20.0	67.4	21.4
Coher Command R+	70.1	28.0	57.2	32.0	85.8	-	70.1	30.9	89.4	34.0	66.8	29.6
Mistral Large	81.7	38.5	<b>78.3</b>	41.2	<b>91.9</b>	-	<b>83.2</b>	40.4	92.9	33.0	71.5	35.6
Mixtral 8x22B MoE	48.9	30.0	50.0	36.8	89.2	-	59.1	35.0	92.1	30.0	74.8	34.1
Llama3 8B	26.6	16.0	20.6	23.7	32.4	-	26.2	21.6	68.8	21.0	54.7	23.6
Llama3 70B	42.8	31.0	61.1	45.5	75.0	-	56.1	41.5	93.2	25.0	63.4	27.7

1712 **Grammar Errors** Figure 23 illustrates parsing errors, predicate parameter length errors, and hallucination  
 1713 errors for both VirtualHome and BEHAVIOR . Specifically: (1) Parsing errors: the VirtualHome example  
 1714 duplicates ‘.’ in an object name, whereas the BEHAVIOR example repeatedly outputs ‘.’. (2) Predicate  
 1715 parameter length errors: the VirtualHome example outputs two parameters for the action *GRAB*, while the  
 1716 BEHAVIOR example outputs three parameters for the state *nextto*. (3) Hallucination errors: the objects *kitchen.1*  
 1717 and *countertop.84* do not exist in the provided environment.

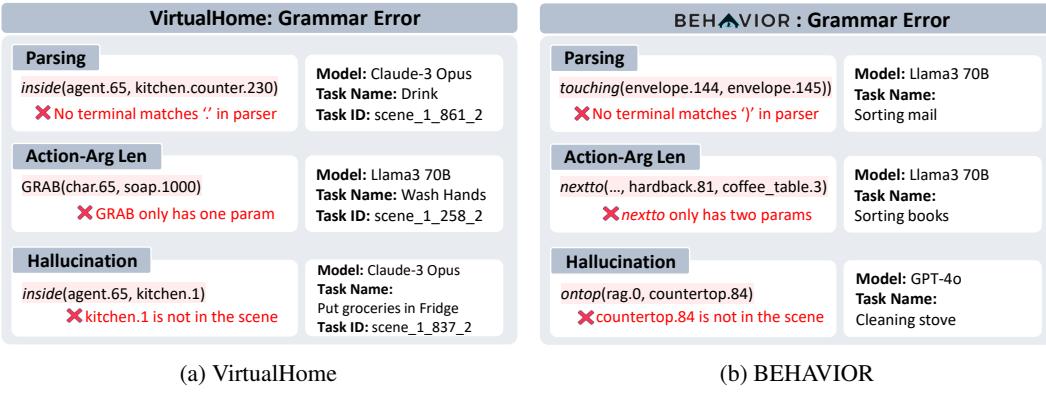


Figure 23: Parsing error examples for subgoal decomposition.

1718 **Runtime Errors** Figure 24 demonstrates four types of runtime errors for both VirtualHome and BEHAVIOR : wrong order, missing step, affordance, and additional step errors. Specifically: (1) Wrong order errors:

1720 In the VirtualHome example, the agent needs to plug in *light.411* but cannot do so because both hands hold  
 1721 *novel.1000* and *spectacles.1001*, whereas they were previously empty. Similarly, in the BEHAVIOR example,  
 1722 the agent drops *soap.0* before cleaning *top\_cabinet.25* despite having held the soap beforehand. (2) Missing  
 1723 step errors: In the VirtualHome scenario, the agent is required to retrieve *food.1000*, but it is inside the closed  
 1724 *freezer.289*. The BEHAVIOR example presents the agent tasked with placing *vidalia\_onion.67* on *countertop.26*,  
 1725 yet the onion is inside a closed container. (3) Affordance errors: For VirtualHome , the agent is asked to plug in  
 1726 *computer.417* and *mouse.413*, but these items lack plugs in VirtualHome . In the BEHAVIOR example, the agent  
 1727 is asked to hold *fridge.97* and slice *carving\_knife.0*. However, the fridge is too large to hold, and the knife cannot  
 1728 be sliced. (4) Additional step errors: In VirtualHome , the agent is instructed to plug in *washing\_machine.1001*,  
 1729 a step already fulfilled initially. Similarly, in BEHAVIOR , the agent is asked to open *carton.0*, which was  
 1730 already open in a previous state.

VirtualHome: Trajectory – Runtime Error					
Wrong Order	Missing Step	Affordance Error	Additional Step		
<b>Model:</b> Claude-3 Opus <b>Task Name:</b> Read book <b>Task ID:</b> scene_1_688_2	<b>Model:</b> Claude-3 Opus <b>Task Name:</b> Put groceries <b>Task ID:</b> scene_1_609_2	<b>Model:</b> Claude-3 Opus <b>Task Name:</b> Work <b>Task ID:</b> scene_1_670_2	<b>Model:</b> Claude-3 Opus <b>Task Name:</b> Wash clothes <b>Task ID:</b> scene_1_27_2		
<i>holds_rh(novel.1000)</i> <i>... holds_lh(spectacles.1001)</i> <i>plugged_in(light.411)</i>	<i>holds_rh(food.1000)</i> <i>nextto(agent.65, freezer.289)</i> <i>...</i>	<i>plugged_in(computer.417)</i> <i>plugged_in(mouse.413)</i> <i>...</i>	<i>next_to(agent.65, washing_machine.1001)</i> <i>plugged_in(washing_machine.1001)</i> <i>...</i>		
<i>✗ Precondition has_a_free_hand(agent.65) = False</i> <i>✓ Historical State has_a_free_hand(agent.65) = True</i>	<i>✗ Precondition not_inside{food.1000) = False</i> <i>✗ Historical State not_inside{food.1000) = False</i>	<i>✗ Precondition computer.417 has_no_plug</i> <i>✗ Precondition mouse.413 has_no_plug</i>	<i>✗ Current State plugged_in(wm.1001) = True</i> <i>✗ Expected State plugged_in(wm.1001) = False</i>		

(a) VirtualHome

BEHAVIOR: Trajectory – Runtime Error					
Wrong Order	Missing Step	Affordance Error	Additional Step		
<b>Model:</b> GPT-4o <b>Task Name:</b> Cleaning kitchen cupboard <i>...</i>	<b>Model:</b> GPT-4o <b>Task Name:</b> Chopping vegetables <i>...</i>	<b>Model:</b> Llama3 70B <b>Task Name:</b> Bottling fruit <i>...</i>	<b>Model:</b> GPT-4o <b>Task Name:</b> Setting up candles <i>...</i>		
<i>holds_lh(soap.0)</i> <i>not_holds_lh(soap.0)</i> <i>not_dusty(top_cabinet.25)</i> <i>...</i>	<i>ontop(vidalia_onion.67, countertop.26)</i> <i>holds_rh(carving_knife.69)</i> <i>...</i>	<i>holds_lh(fridge.97)</i> <i>not_inside(strawberry.0, fridge.97)</i> <i>sliced(carving_knife.0)</i>	<i>open(carton.0)</i> <i>holds_rh(candle.0)</i> <i>ontop(candle.0, coffee_table.33)</i> <i>...</i>		
<i>✗ Precondition holding(soap.0) = False</i> <i>✓ Historical State holding(soap.0) = True</i>	<i>✗ Precondition not_inside{vidalia_onion.67) = False</i> <i>✗ Historical State not_inside{vidalia_onion.67) = False</i>	<i>✗ Precondition fridge.97 too_big_to_grasp</i> <i>✗ Precondition carving_knife.0 not_sliceable</i>	<i>✗ Current State open(carton.0) = True</i> <i>✗ Expected State open(carton.0) = False</i>		

(b) BEHAVIOR

Figure 24: Trajectory runtime error examples for subgoal decomposition.

1731 **Goal Satisfaction Errors** Figure 25 depicts missing state, missing relation errors for both VirtualHome  
 1732 and BEHAVIOR , and includes an example of a missing goal action error for VirtualHome . Specifically: (1)  
 1733 Missing state errors: In the VirtualHome example, the generated subgoal plan fails to turn on *light.1002*, which  
 1734 is necessary for achieving the final goal. Similarly, in the BEHAVIOR example, the subgoal plan fails to open  
 1735 *window.76* and *window.81*, which are also required for the final goal. (2) In VirtualHome , after execution, the  
 1736 agent does not place *clothes\_pants.1001* on top of *washing\_machine.1000*, as required by the final goal. In the  
 1737 BEHAVIOR example, LLMs incorrectly instruct the agent to first put *plywood.78* next to *plywood.79*, and then  
 1738 put *plywood.79* next to *plywood.80*. Once the second state is achieved, the state *nextto(plywood.78, plywood.79)*  
 1739 is no longer satisfied because *plywood.79* has been moved. (3) Missing goal action error: In the VirtualHome  
 1740 example, the agent is required to perform the action *DRINK*, but this action is missing in the generated subgoal  
 1741 plan.

## Takeaways of Subgoal Decomposition

- (1) **Top-performing Models:** GPT-4o demonstrates superior performance in both VirtualHome and BEHAVIOR simulator compared with other SOTA LLMs. Gemini 1.5 Flash also demonstrates highest performance in VirtualHome with success rate close to 90%.
- (2) **Grammar and Runtime Errors:** State-of-the-art models generally avoid grammar errors but can hallucinate actions and objects. The most common runtime errors in VirtualHome are additional steps, whereas most of runtime errors in BEHAVIOR are missing steps.
- (3) **Goal Satisfaction:** Stronger LLMs like GPT-4o show higher accuracy in action goal success rates in VirtualHome compared to weaker models like Llama3 8B. However, achieving state and relation goals in BEHAVIOR is challenging due to more complex task representations and stricter precondition checks.
- (4) **Impact of Task Complexity In Different Simulators:** Overall LLM performance is lower in BEHAVIOR compared with VirtualHome due to complex task representations involving quantifiers like *forall* and *forpairs*, which articulate complex temporal and spatial requirements.

1742

### 1743 H.3 Action Sequencing

1744 The full results for the trajectory evaluation of action sequencing are presented in Table 14. Detailed results  
 1745 for goal satisfaction are shown in Table 13. The results from two simulators are provided separately, with *V*  
 1746 representing VirtualHome and *B* denoting BEHAVIOR .

1747 **(1) Executable and Goal Success Rate.** For VirtualHome , Mixtral Large achieved the best performance in  
 1748 both execution and goal success rates, which is surprising as it outperforms some LLMs that're considered to be  
 1749 more powerful such as Cluade-3 Opus, GPT-4o and Gemini-1.5 Pro. And the same case applies to Cluade-3  
 1750 Sonnet, which outperforms Cluade-3 Opus in terms of execution success rate. We hypothesize that this is due to  
 1751 the fact that most tasks in VirtualHome require only 3-5 actions, but may have a large number of objects in the  
 1752 scene. Therefore, it is important to capture the key information and focus only on task-related objects. In such  
 1753 scenarios, some LLMs may "overthink" and perform actions on irrelevant objects, while LLMs like Mixtral can  
 1754 think more straightforwardly to capture the key information.

1755 For BEHAVIOR , the situation is the opposite: most tasks require long sequences of execution, with an average  
 1756 of 14.6 steps as shown in Figure 55. Additionally, the preconditions for the actions in BEHAVIOR are more  
 1757 complex and strict compared to VirtualHome , while the number of interactable objects in the scene is relatively  
 1758 small. This requires more long-term planning and commonsense reasoning abilities. Thus, LLMs that "think  
 1759 more" win, explaining why the closed-source and powerful LLMs such as Claude-3 Opus achieved the best  
 1760 performance.

1761 In general, the performance of LLMs in VirtualHome is significantly better than in BEHAVIOR , which is  
 1762 reasonable considering the task complexity. However, it is noteworthy that even the best LLMs only completed  
 1763 half of the tasks in BEHAVIOR which indicates there's still a long way to go before LLMs can be effectively  
 1764 applied to embodied agents tasks.

1765 **(2) Grammar Errors.** The occurrence of grammar errors for different LLMs generally aligns with our overall  
 1766 understanding of their capabilities. That is, the larger the model, the fewer the grammar errors; newer editions  
 1767 tend to perform better than older ones; and LLMs trained with high-quality data excel. For instance, Claude-  
 1768 3 Opus, GPT-4o, and Mistral Large made fewer grammar errors compared to their other family members,  
 1769 indicating advancements in their language understanding and generation capabilities. This trend holds true for  
 1770 both BEHAVIOR and VirtualHome tasks.

VirtualHome: Goal Satisfaction Error			BEHAVIOR: Goal Satisfaction Error		
<b>Missing State</b> Model: Claude-3 Opus Task Name: Turn on light Task ID: scene_1_160_1	<b>Missing Relation</b> Model: GPT-4o Task Name: Wash clothes Task ID: scene_1_850_1	<b>Missing Goal Action</b> Model: Llama3 8B Task Name: Drink Task ID: scene_1_1064_1	<b>Missing State</b> Model: Cohere Command R Task Name: Locking every window	<b>Missing Relation</b> Model: GPT-4o Task Name: Laying wood floors	
<b>Goal</b> <code>on[light.411] and on[light.1002]</code>	<b>Goal</b> <code>ontop(clothes_pants.1001, washing_machine.1000)</code>	<b>Goal</b> <code>DRINK</code>	<b>Goal</b> <code>not open(window.76) not open(window.81)</code>	<b>Goal</b> <code>nextto(plywood.78, plywood.79) and nextto(plywood.79, plywood.80)</code>	
<b>LLM Output</b> <code>plugged_in(light.411) clean(light.411) on(light.411)</code>	<b>LLM Output</b> <code>open(washing_machine.1000) ontop(laundry_detergent.1002, washing_machine.1000)</code>	<b>LLM Output</b> <code>NEXT_TO(agent.65, water_glass.1000) HOLDS_RH(agent.65, water_glass.1000)</code>	<b>LLM Output</b> <code>open(window.76) open(window.81)</code>	<b>LLM Output</b> <code>nextto(plywood.78, plywood.79) nextto(plywood.79, plywood.80) nextto(plywood.80, plywood.81)</code>	<b>Error Info: Relation Unsatisfied</b> <code>nextto(plywood.78, plywood.79)</code>
<b>Error Info: State Unsatisfied</b> <code>on(light.1002)</code>	<b>Error Info: Missing Final Relation</b> <code>ontop(clothes_pants.1001, washing_machine.1000)</code>	<b>Error Info: Action Unsatisfied</b> <code>DRINK</code>	<b>Error Info: State Unsatisfied</b> <code>not open(window.76) not open(window.81)</code>	<b>Error Info: Relation Unsatisfied</b> <code>nextto(plywood.78, plywood.79)</code>	<b>Error Info: Missing Final State</b> <code>not open(window.76) not open(window.81)</code>

(a) VirtualHome

(b) BEHAVIOR

Figure 25: Goal satisfaction error examples for subgoal decompositions

1771 We studied the cases where LLMs made grammar errors, as shown in Figure 26. Common parsing errors include  
1772 unfinished actions or action sequences, often due to LLMs mistakenly stopping generation early or producing  
1773 overly long and incorrect action sequences that exceed the context window. Additionally, LLMs sometimes  
1774 confuse the length of arguments. For example, although GPT-4o made very few grammar errors, it still confuses  
1775 the parameters required for *rinse* and *wash* in VirtualHome , leading to action-arg length errors. Interestingly,  
1776 *rinse* seems particularly challenging for GPT-4o, as it also made another hallucination error involving *rinse*.

1777 We consider this metric a reflection of the LLMs’ instruction-following and trustworthiness abilities. LLMs  
1778 that can better follow instructions are likely to produce more grammatically correct outputs, which is crucial for  
1779 applications requiring high levels of precision and reliability. This underscores the importance of continuous  
1780 improvements in model size, training methodologies, and data quality to enhance the reliability and accuracy of  
1781 language LLMs.

1782 **(3) Runtime Errors.** The runtime error rate in VirtualHome is lower than BEHAVIOR , the reason of which has  
1783 been previously discussed: 1) the action sequences required to accomplish tasks in BEHAVIOR are longer, and  
1784 2) the preconditions for executing an action are more complex and strict. For both simulators, the most frequent  
1785 runtime errors are missing steps, as unsatisfied preconditions are the primary cause of failed action execution.

1786 Figure 27 shows examples of errors made by LLMs. The reasons for making a *missing step* error are mainly due  
1787 to the lack of common sense reasoning ability to satisfy the preconditions for a desired action, e.g., *getting to*  
1788 *the sink*, *holding a soap* before *washing hands*, or *soaking the brush* before *cleaning a stain*. *Wrong temporal*  
1789 *order* errors often arise from deficiencies in planning abilities, such as mistakenly *drinking* after *putting a cup*  
1790 *back* or *grasping a tomato* before *slicing it*, demonstrating that LLMs still lack experience or common sense  
1791 with daily life activities. Affordance errors occur due to a lack of understanding about properties of objects, like  
1792 attempting to *type on a mouse* or not *assuming a strawberry will become different parts after slicing it*. For  
1793 *additional step* errors, it may be an issue with in-context memory, where LLMs forget they have already opened  
1794 a cabinet and attempt to open it again, or it could be related to reasoning about the common situation of the  
1795 initial scene, like *the agent would be standing at the beginning*, therefore no need to *instruct it to stand up again*.  
1796 In sum, the ability of commonsense reasoning is crucial for successfully predicting action sequences since our  
1797 prompts do not provide detailed information about the environment and action instructions. LLMs must follow  
1798 human conventions to reasonably guess the preconditions and post-effects of the actions.

1799 For the runtime error metrics, none of the LLMs significantly outperformed the others. While some commonly  
1800 acknowledged powerful LLMs, such as GPT-4-turbo, Gemini 1.5 Pro, and Claude-3 Opus, performed well in  
1801 certain areas, they still struggled with other aspects of the metrics.

1802 **(4) Goal Satisfaction.** In general, in VirtualHome , LLMs perform better at satisfying state goals than relation  
1803 goals, while in BEHAVIOR , it is vice versa. The reason for this difference is that in VirtualHome , there  
1804 are more objects in the scene, and objects can have more complex relations (like triple relations) compared to  
1805 BEHAVIOR , where only binary relations exist. In BEHAVIOR , the preconditions for state actions are more  
1806 complex, as shown in Tables 25 and 24. Therefore, it is easier for LLMs to achieve state goals in VirtualHome  
1807 and relation goals in BEHAVIOR .

1808 However, performance varies across different LLMs. In both simulators, GPT-4-turbo achieves the highest state  
1809 goal satisfaction rates, which are higher than its relation goal satisfaction rates. Additionally, Claude-3 Sonnet  
1810 and Mixtral Large performed exceptionally well in relation goals.

1811 Figure 28 shows examples of errors that prevent LLMs from satisfying goals in a task. Generally, there are two cases where LLMs fail to satisfy a goal: 1) Missed goals: forgetting to generate actions to achieve the goal, e.g., forgetting to *TOGGLE\_ON(laptop.1000)* to satisfy *on(laptop.1000)*, which could be due to deficiencies in instruction-following or issues with in-context memory. 2) Wrong actions: mistakenly corresponding goals to incorrect actions, e.g., corresponding *onfloors(plywood\_78)* to *RIGHT\_PLACE\_NEXTTO(room\_floor\_living\_room\_0)* instead of *RIGHT\_PLACE\_FLOOR(room\_floor\_living\_room\_0)*, which could be a lack of reasoning ability.

Table 13: Goal satisfaction rates (%) for *action sequencing*. Full results.

Model	State Goal		Relation Goal		Action Goal		Total	
	V	B	V	B	V	B	V	B
Claude-3 Haiku	59.4	27.0	42.8	38.7	87.8	-	61.4	35.5
Claude-3 Sonnet	79.5	41.0	67.2	<b>59.8</b>	85.1	-	77.2	<b>54.6</b>
Claude-3 Opus	63.7	45.0	71.1	53.0	77.0	-	69.1	50.8
Gemini 1.0 Pro	52.5	28.0	32.8	32.0	77.0	-	52.6	30.9
Gemini 1.5 Flash	79.5	34.0	63.3	50.0	88.5	-	76.9	45.6
Gemini 1.5 Pro	81.7	41.0	76.7	43.2	89.2	-	82.0	42.6
GPT-3.5-turbo	29.1	20.0	15.6	22.6	64.2	-	33.7	21.9
GPT-4-turbo	74.8	39.0	72.2	39.5	89.9	-	77.7	39.3
GPT-4o	<b>83.1</b>	<b>49.0</b>	71.1	45.5	89.9	-	81.2	46.5
Cohere Command R	18.4	20.0	31.1	25.9	48.0	-	29.4	24.3
Cohere Command R+	70.1	28.0	57.2	32.0	85.8	-	70.1	30.9
Mistral Large	81.7	38.5	<b>78.3</b>	41.2	<b>91.9</b>	-	<b>83.2</b>	40.4
Mixtral 8x22B MoE	48.9	30.0	50.0	36.8	89.2	-	59.1	35.0
Llama3 8B	26.6	16.0	20.6	23.7	32.4	-	26.2	21.6
Llama3 70B	42.8	31.0	61.1	45.5	75.0	-	56.1	41.5

### Takeaways of Action Sequencing

1. Larger models did not always outperform smaller ones, especially for shorter and simpler tasks.
2. No LLMs (including powerful ones like GPT-4-turbo, Gemini 1.5 Pro, and Claude-3 Opus) consistently outperformed others across all runtime error metrics.
3. Generally, larger model sizes and higher quality training data correlated with lower grammar error rates.
4. Commonsense reasoning about object properties, agent states, and preconditions for actions remains a significant challenge for current LLMs, particularly when given limited context.
5. Errors in goal satisfaction can be attributed to: 1) Missed goals: Failing to generate actions to achieve the goal. 2) Wrong actions: Incorrectly mapping goals to inappropriate actions.

1818

Table 14: Trajectory evaluation results (%) for *action sequencing*. Full results.

Model	Goal Evaluation				Trajectory Evaluation							
	Goal SR		Execution SR		Grammar Error (↓)				Runtime Error (↓)			
	V	B	V	B	Parsing	Hallucination	Action-Arg Num	Wrong Order	Missing Step	Affordance	Additional Step	V
Claude-3 Haiku	43.6	26.0	51.5	32.0	<b>0.0</b>	<b>0.0</b>	4.9	6.0	0.3	<b>0.0</b>	7.0	42.0
Claude-3 Sonnet	65.2	44.0	68.9	57.0	<b>0.0</b>	<b>0.0</b>	5.6	1.0	0.7	<b>0.0</b>	0.7	11.0
Claude-3 Opus	65.9	<b>51.0</b>	64.9	<b>59.0</b>	<b>0.0</b>	<b>0.0</b>	14.1	<b>0.0</b>	<b>0.0</b>	1.3	3.0	19.0
Gemini 1.0 Pro	33.1	27.0	36.7	32.0	0.7	7.0	9.2	3.0	10.5	6.0	0.3	13.0
Gemini 1.5 Flash	61.0	40.0	65.9	52.0	<b>0.0</b>	<b>0.0</b>	2.0	<b>0.0</b>	0.3	<b>0.0</b>	5.0	30.8
Gemini 1.5 Pro	75.1	42.0	82.0	54.0	0.3	<b>0.0</b>	1.6	<b>0.0</b>	0.3	<b>0.0</b>	6.0	14.8
GPT-3.5-turbo	25.9	16.0	40.7	20.0	<b>0.0</b>	4.0	4.3	7.0	17.7	23.0	<b>0.0</b>	<b>1.0</b>
GPT-4-turbo	60.7	38.0	64.6	45.0	<b>0.0</b>	<b>0.0</b>	1.6	<b>0.1</b>	<b>0.0</b>	7.0	32.1	47.0
GPT-4o	70.2	47.0	71.8	53.0	<b>0.0</b>	<b>0.0</b>	<b>1.3</b>	1.0	0.7	<b>0.0</b>	9.0	25.3
Cohere Command R	15.7	16.0	19.7	19.0	2.0	5.0	37.1	13.0	23.3	<b>0.0</b>	0.3	8.0
Cohere Command R+	54.8	27.0	61.0	35.0	<b>0.0</b>	<b>0.0</b>	5.9	1.0	2.6	15.0	<b>0.0</b>	10.0
Mistral Large	<b>78.4</b>	33.0	<b>83.9</b>	50.0	<b>0.0</b>	<b>0.0</b>	3.3	<b>0.0</b>	0.3	<b>0.0</b>	8.0	<b>12.5</b>
Mixtral 8x22B MoE	46.2	30.0	50.2	40.0	<b>0.0</b>	3.0	13.1	6.0	0.7	<b>0.0</b>	10.0	34.8
Llama3 8B	22.3	10.0	22.3	16.0	<b>0.0</b>	<b>0.0</b>	43.6	15.0	5.6	<b>0.0</b>	6.0	28.5
Llama3 70B	54.4	34.0	56.1	42.0	<b>0.0</b>	<b>0.0</b>	23.3	2.0	1.0	<b>0.0</b>	0.7	15.0

1819 

## H.4 Transition Modeling

1820 

### H.4.1 Logic Form Accuracy

1821 **VirtualHome** The results from VirtualHome (Table 15) illustrate the varied performance of models across five  
1822 distinct categories, demonstrating their ability to predict complex logic and predicates for each action. The  
1823 model Claude-3 Opus showcased superior performance in terms of all metrics in the object states category,  
1824 achieving an F1 score of 63.0%, highlighting its effectiveness in interpreting complex object state transitions.  
1825 Notably, Gemini 1.5 Pro performed exceptionally well in object orientation, achieving the highest F1 score of  
1826 90.9%, which suggests its strong capability in understanding object orientations like 'facing'. Gemini 1.5 Pro  
1827 also excels in the object affordance category on the precision, recall and f1 score with an F1 score of 77.7%,  
1828 demonstrating its well understanding of required object properties for different actions. However, across all  
1829 models, the non-spatial relations category displayed generally low scores, with Gemini 1.5 Flash performing  
1830 best at a modest F1 of 7.9%. This basically results from the complex logic and corner cases involved in some

VirtualHome: Grammar Error		BEHAVIOR: Grammar Error	
<b>Parsing</b>	<b>Model:</b> Gemini 1.5 Pro <b>Task Name:</b> Wash hands <b>Task ID:</b> scene_1_657_1  ✖️ Unfinished action sequence with '.', at the end	<b>Model:</b> GPT-3.5 <b>Task Name:</b> Setting up candles  ✖️ Unfinished action, no object	
<b>Action-Arg Len</b>	<b>Model:</b> GPT-4o <b>Task Name:</b> Wash Hands <b>Task ID:</b> scene_1_258_2  ✖️ RINSE and WASH should have one param	<b>Model:</b> GPT-3.5 <b>Task Name:</b> Brushing lint off clothing  ✖️ LEFT_PLACE_ONTOP only has one param	
<b>Hallucination</b>	<b>Model:</b> GPT-4o <b>Task Name:</b> Wash hands <b>Task ID:</b> scene_1_580_1  ✖️ faces not in the scene	<b>Model:</b> Claude-3 Sonnet <b>Task Name:</b> Setting mousetraps  ✖️ LEFT_PLACE_ONFLOOR is not a valid action	

Figure 26: Grammar error examples for *action sequencing*.

VirtualHome: Trajectory – Runtime Error			
Wrong Order	Missing Step	Affordance Error	Additional Step
<b>Model:</b> Gemini 1.5 Flash <b>Task Name:</b> Drink <b>Task ID:</b> scene_1_171_2  ... PUTBACK(cup,100 0,sink,231) DRINK(cup,1000) ...  ✖️ Precondition holds(cup,1000) = False ✓ Historical State holds(cup,1000) = True	<b>Model:</b> Gemini 1.5 Flash <b>Task Name:</b> Wash hands <b>Task ID:</b> scene_1_813_2  WALK(bathroom,1) RINSE(hands_both .1000) ...  ✖️ Precondition next_to(sink,42) = False holds(soap,100) = False ✖️ Historical State next_to(sink,42) = False holds(soap,100) = False	<b>Model:</b> Mixtral 8x22b MOE .319) FIND(mouse,413) TYPE(mouse,413) ...  ✖️ Affordance mouse,413 can't be typed	<b>Model:</b> Mistral Large ... STANDUP()  ✖️ Current State stand_up(character,45) = True ↓ Expected State stand_up(character,45) = False
<b>Model:</b> GPT-4o <b>Task Name:</b> Chopping Vegetables <b>Task ID:</b> scene_1_670_2  RIGHT_GRASP(car ving_knife_69) LEFT_GRASP(toma to_61) SLICE(tomato_61)  ✖️ Precondition not_in_hand(tomato_61)= False ✓ Historical State not_in_hand(tomato_61)= True	<b>Model:</b> GPT-4o <b>Task Name:</b> Cleaning bathtubs <b>Task ID:</b> scene_1_279_2  stained(bathtub_35) RIGHT_GRASP(scr ub_brush_0) CLEAN(bathtub_35)  ✖️ Precondition soaked(scrub_brush_0) = False ✖️ Historical State soaked(scrub_brush_0) = False	<b>Model:</b> Claude-3 Sonnet <b>Task Name:</b> Bottling fruit <b>Task ID:</b> scene_1_93_1  sliced(strawberry_ 0) RIGHT_TRANSFER_ CONTENTS_INSIDE(stawberry_0)  ✖️ Affordance strawberry_0 is sliced and not interactable. Should interact with strawberry_0_part0 and strawberry_0_part1	<b>Model:</b> Claude-3 Opus OPEN(top_cabinet_27 ... OPEN(top_cabinet_27)  ✖️ Current State open(top_cabinet_27) = True ↓ Expected State open(top_cabinet_27) = False

Figure 27: Trajectory runtime error examples for *action sequencing*.

VirtualHome: Goal Satisfaction Error		
Missing State	Missing Relation	Missing Goal Action
<p><b>Model:</b> Gemini 1.5 Pro <b>Task Name:</b> Write an email <b>Task ID:</b> scene_1_87_2</p> <p><b>Goal</b> on(laptop,1000) and TYPE(laptop,1000)</p> <p><b>LLM Output</b> WALK(laptop) TYPE(laptop)</p> <p><b>Error Info: State Unsatisfied</b> ✖ Missing Final State on(laptop,1000)</p>	<p><b>Model:</b> Gemini 1.5 Pro <b>Task Name:</b> Wash teeth <b>Task ID:</b> scene_1_312_2</p> <p><b>Goal</b> holds_lh(toothbrush,1001)</p> <p><b>LLM Output</b> WALK(toothbrush) GRAB(toothbrush)</p> <p><b>Error Info: Relation Unsatisfied</b> ✖ Missing Final Relation holds_lh(toothbrush,1001)</p>	<p><b>Model:</b> Llama 3 8b <b>Task Name:</b> Pet cat <b>Task ID:</b> scene_1_137_2</p> <p><b>Goal</b> TOUCH(cat,1000)</p> <p><b>LLM Output</b> WALK(cat)</p> <p><b>Error Info: Action Unsatisfied</b> ✖ Missing Goal Action TOUCH(cat,1000)</p>
BEHAVIOR: Goal Satisfaction Error		
<p><b>Missing State</b></p> <p><b>Model:</b> GPT-4o <b>Task Name:</b> Washing floor</p> <p><b>Goal</b> ["nor", "or", "dusty", "room_floor_bathroom_0", "stained", "room_floor_bathroom_0"]</p> <p><b>LLM Output</b> RIGHT_GRASP(bath_towel_0) CLEAN(room_floor_bathroom_0) ...</p> <p><b>Error Info: State Unsatisfied</b> ✖ Missing Final State (not stained room_floor_bathroom_0)</p>	<p><b>Missing Relation</b></p> <p><b>Model:</b> GPT-4o <b>Task Name:</b> Laying wood floors</p> <p><b>Goal</b> ["forall", "plywood.n.01", ":", "plywood.n.01", "onfloor", "plywood.n.01", "room_floor_living_room_0"]</p> <p><b>LLM Output</b> RIGHT_GRASP(plywood_78) RIGHT_PLACE_NEXTTO(room_floor_living_room_0)</p> <p><b>Error Info: Relation Unsatisfied</b> ✖ Missing Final Relation onfloor(plywood_78)</p>	<p><b>Missing Relation</b></p> <p><b>Model:</b> GPT-4o <b>Task Name:</b> Thawing frozen food</p> <p><b>Goal</b> ["exists", "container_date_0", "fish_0"]</p> <p><b>LLM Output</b> RIGHT_GRASP(container_date_0) RIGHT_PLACE_NEXTTO(fish_0) RIGHT_GRASP(fish_0) RIGHT_PLACE_NEXTTO(link_37)</p> <p><b>Error Info: Relation Unsatisfied</b> ✖ Missing Final Relation (notco container_date_0 fish_0)</p>

Figure 28: Goal satisfaction error examples for *action sequencing*.

Table 15: Full results of logic form accuracy for *transition modeling* in VirtualHome

Model	Object States			Object Orientation			Object Affordance			Spatial Relations			Non-Spatial Relations		
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
Claude-3 Haiku	76.0	40.1	52.5	19.0	34.4	24.4	67.8	73.9	70.7	37.7	38.7	38.2	2.0	1.5	1.7
Claude-3 Opus	<b>87.4</b>	<b>49.2</b>	<b>63.0</b>	46.3	<b>96.9</b>	62.6	76.8	74.3	75.5	37.6	39.9	38.7	10.4	<b>5.2</b>	7.0
Claude-3 Sonnet	76.6	37.4	50.3	48.1	78.1	59.5	60.7	74.3	66.8	32.3	39.9	35.7	6.2	4.1	4.9
Cohere Command R	18.0	6.8	9.9	38.7	90.6	54.2	40.2	23.0	29.2	12.6	6.7	8.8	3.3	0.9	1.4
Cohere Command R+	44.9	19.0	26.3	34.6	68.8	45.9	51.0	62.1	56.0	30.1	34.8	32.4	7.6	3.1	4.4
Gemini 1.0 Pro	68.4	12.3	20.4	16.3	62.5	27.9	55.3	20.1	29.6	45.0	16.5	24.3	7.7	2.5	3.8
Gemini 1.5 Flash	82.3	37.6	51.6	2.0	3.1	2.5	54.4	74.7	62.9	<b>47.4</b>	<b>42.9</b>	<b>45.0</b>	<b>16.3</b>	<b>5.2</b>	<b>7.9</b>
Gemini 1.5 Pro	45.3	11.9	18.8	<b>88.2</b>	93.8	<b>90.9</b>	<b>79.9</b>	<b>75.5</b>	<b>77.7</b>	42.2	35.8	38.7	15.5	<b>5.2</b>	7.8
GPT-3.5-turbo	63.5	21.9	32.5	11.4	15.6	13.2	57.2	53.1	54.9	35.2	21.7	26.8	1.7	0.3	0.6
GPT-4-turbo	79.3	44.2	56.7	10.1	31.3	15.3	65.9	71.0	68.4	31.8	34.2	32.9	3.8	1.0	1.6
GPT-4o	80.2	41.5	54.6	48.0	59.4	52.8	76.2	73.7	74.9	40.8	40.7	40.8	14.8	5.1	7.5
Llama3 8b	30.8	13.7	18.9	0.0	0.0	0.0	1.6	3.2	2.1	15.5	18.2	16.8	0.0	0.0	0.0
Llama3 70b	63.5	21.9	32.5	49.0	66.3	56.6	65.0	50.0	57.0	27.0	27.0	27.0	5.0	2.0	3.0
Mistral Large	30.0	8.0	13.0	48.0	88.0	62.0	72.0	29.0	41.0	35.0	18.0	24.0	3.0	1.0	1.0
Mixtral 8x22B MoE	72.0	33.0	45.0	43.0	83.0	57.0	64.0	74.0	69.0	40.0	38.0	39.0	12.0	4.0	6.0

Table 16: Full results of logic form accuracy for *transition modeling* in BEHAVIOR

Model	Object States			Spatial Relations			Non-Spatial Relations		
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
Claude-3 Haiku	83.9	90.0	86.8	90.0	89.9	89.9	91.8	88.6	90.2
Claude-3 Opus	<b>91.6</b>	<b>99.3</b>	<b>95.3</b>	91.6	94.5	93.1	93.1	<b>90.0</b>	<b>91.6</b>
Claude-3 Sonnet	90.7	97.8	94.1	92.8	94.3	93.5	92.7	85.7	89.1
Cohere Command R	82.4	59.9	69.4	58.7	43.2	49.8	69.1	57.4	62.7
Cohere Command R+	58.9	58.2	58.6	81.0	67.7	73.8	70.3	83.2	76.2
Gemini 1.0 Pro	79.7	87.7	83.5	86.2	93.2	89.6	79.7	79.1	79.4
Gemini 1.5 Flash	69.8	75.5	72.6	89.9	94.0	91.9	84.6	88.9	86.7
Gemini 1.5 Pro	90.0	91.7	90.8	91.6	86.4	88.9	91.0	83.6	87.1
GPT-3.5-turbo	74.8	68.8	71.7	71.6	84.9	77.7	87.9	66.9	75.9
GPT-4-turbo	78.8	75.6	77.2	84.4	65.0	73.5	<b>94.5</b>	72.8	82.2
GPT-4o	89.7	96.9	93.2	<b>94.2</b>	<b>95.3</b>	<b>94.8</b>	93.4	89.9	<b>91.6</b>
Llama3 8b	61.9	47.6	53.8	71.3	54.8	62.1	87.8	70.3	78.0
Llama3 70b	90.7	95.2	92.9	76.8	75.7	76.2	90.8	65.6	76.2
Mistral Large	84.1	85.4	84.7	79.8	70.7	75.0	89.7	80.2	84.7
Mixtral 8x22B MoE	85.2	93.8	89.3	89.3	91.3	90.3	83.7	86.3	85.0

non-spatial actions. For example, when predicting action 'grab', few models consider 'not in a closed container' or 'both hands of the robot is holding something' as a precondition, and including the logic 'when the robots hold something with one hand, hold the object with the other hand' in the effect is also too challenging for LLMs. This could point to a common challenge in modeling complex relationships with real-world scenario concerns, requiring perhaps a more nuanced understanding or a different approach in training.

Table 17: Full results of planner success rate for *transition modeling (%)*

Model	Object States		Object Orientation		Object Affordance		Spatial Relations		Non-Spatial Relations	
	V	B	V	B	V	B	V	B	V	B
Claude-3 Haiku	13.5	68.9	3.6	-	19.8	-	46.9	80.9	73.0	93.4
Claude-3 Opus	63.5	<b>95.6</b>	71.4	-	58.7	-	64.8	<b>96.8</b>	55.4	<b>98.4</b>
Claude-3 Sonnet	11.2	77.8	3.6	-	10.8	-	20.0	76.6	13.5	80.3
Cohere Command R	44.6	37.8	82.1	-	40.1	-	62.6	33.0	58.3	29.5
Cohere Command R+	36.5	46.7	46.4	-	35.3	-	40.7	30.9	31.1	29.5
Gemini 1.0 Pro	10.7	51.1	0.0	-	10.2	-	14.5	36.2	2.7	31.2
Gemini 1.5 Flash	34.8	77.8	7.1	-	46.7	-	61.4	74.5	60.8	73.8
Gemini 1.5 Pro	<b>94.4</b>	77.8	<b>89.3</b>	-	<b>95.8</b>	-	<b>89.0</b>	64.9	<b>83.8</b>	62.3
GPT-3.5-turbo	1.1	42.2	25.0	-	1.2	-	0.0	41.5	0.0	32.8
GPT-4-turbo	51.7	88.9	50.0	-	47.9	-	67.6	83.0	64.9	82.0
GPT-4o	71.9	91.1	78.6	-	63.5	-	66.9	90.4	68.9	88.5
Llama3 8b	27.0	68.9	0.0	-	26.4	-	37.9	43.6	31.1	36.1
Llama3 70b	10.1	17.8	3.6	-	6.6	-	15.2	10.6	18.9	9.9
Mistral Large	15.7	53.3	7.1	-	14.4	-	17.9	27.7	8.1	23.0
Mixtral 8x22B MoE	36.5	68.9	50.0	-	28.1	-	44.1	48.9	43.2	44.3

1836 **BEHAVIOR** In the BEHAVIOR environment (Table 16), the model Claude-3 Opus again shows outstanding  
 1837 performance, particularly in object states with an F1 score of 95.3% and in non-spatial relations with an F1  
 1838 of 91.6%. These results underscore its exceptional ability to handle both static state transitions of objects and  
 1839 non-spatial relations between the robot and the objects. The GPT-4o model displayed remarkable precision in  
 1840 Spatial Relations, achieving an F1 score of 94.8%, the highest among all models. This suggests that GPT-4o  
 1841 is particularly adept at modeling spatial inter-object dynamics. Besides, the overall high score in BEHAVIOR  
 1842 results from the small state, action, and predicate space of BEHAVIOR and stronger in-context learning examples  
 1843 we use in BEHAVIOR . The difference between VirtualHome and BEHAVIOR poses challenges for LLM in  
 1844 different environments and test their robustness.

1845 **Overall Conclusions** Therefore, the consistent high performance of Claude-3 models across both environments  
 1846 suggests that these models have potentially benefited from training regimens or architectural features that enhance  
 1847 their understanding of object-oriented and relational aspects of embodied environments. There are also some  
 1848 divergence between the results of two environments. In spatial relation category, Gemini 1.5 Flash performs  
 1849 better in VirtualHome while GPT-4o performs better in BEHAVIOR . The performance disparity between  
 1850 different environments also suggests that model training and optimization might need more customization  
 1851 towards specific types of embodied interactions. Overall, high-performing models in certain categories, such as  
 1852 Gemini 1.5 Pro in Object Orientation and Claude-3 Opus in Object States, highlight areas where specific models  
 1853 excel and can be leveraged for tasks requiring high accuracy in those domains. On the other hand, models like  
 1854 Llama3 and Cohere Command R+ showed lower performance, which may indicate limitations in their training  
 1855 datasets or model architectures for these specific task requirements. Additionally, the generally lower scores in  
 1856 non-spatial Relations across most models suggest a gap in current modeling capabilities, offering a potential  
 1857 area for future research and model improvement. Furthermore,

#### 1858 H.4.2 Planner Success Rate Analysis

1859 This assessment highlights how well models predict feasible transitions and achieve objectives from initial states.

1860 **Model Performance Highlights:** Gemini 1.5 Pro and Claude-3 Opus standout, demonstrating robust success  
 1861 rates across various categories. Gemini 1.5 Pro’s performance in VirtualHome , especially in Object Affordance  
 1862 (95.8% success rate), reflects its precise handling of object functionalities. In contrast, Claude-3 Opus excels  
 1863 in BEHAVIOR, achieving up to 98.4% in Non-Spatial Relations, indicating strong predictive capabilities for  
 1864 complex relational dynamics without spatial constraints.

1865 **Spatial vs. Non-Spatial Relations:** The data reveals a notable disparity in performance between spatial and  
 1866 non-spatial Relations, with non-spatial Relations generally showing higher success rates. This suggests that even  
 1867 though LLM may not cover corner cases in non-spatial relation actions, models can still simplify it and get away  
 1868 with it compared with spatial relation actions, an insight that could direct future training enhancements.

1869 **Challenges and Strategic Insights:** Despite successes, persistent challenges in categories like object orientation,  
 1870 where success rates are generally lower, highlight potential gaps in model training or architectural design. The  
 1871 variability in success rates across models points to significant differences in how they interpret and execute tasks  
 1872 within these dynamic environments. Even high-performing models struggle with precise spatial interactions,  
 1873 a critical area for tasks requiring detailed physical interactions. For example, GPT-4o struggles at non-spatial  
 1874 relation tasks and Claude-3 Opus is not doing well at object states and object affordance tasks.

1875 **Future Directions:** The findings suggest valuable pathways for model improvements, particularly in enhancing  
 1876 relation understanding in physical task execution. Future research might focus on diversifying training scenarios

Missing Predicates			
Missing Object State	Missing Non-spatial Relation	Missing Object Affordance	Missing Spatial Relation
<b>Ground Truth</b> :action open :precondition (and (can_open ?obj)) (closed ?obj) (next_to ?char ?obj)) (not (on ?obj)))	<b>Ground Truth</b> :action put_on :precondition (and (next_to ?char ?obj2) (or (holds_lh ?char ?obj1) (holds_rm ?char ?obj1)))	<b>Ground Truth</b> :action lie :precondition (and (lieable ?obj) (next_to ?char ?obj)) (not (lying ?char)))	<b>Ground Truth</b> :action walk_int :effect (and (inside ?char ?room)) (forall (?far_obj - object) (when (not (not (inside_room ?far_obj ?room)))) (not (next_to ?char ?far_obj))))
<b>LLM Output</b> :action switch_on :precondition (and (can_open ?obj)) (closed ?obj) (next_to ?char ?obj))	<b>LLM Output</b> :action put_on :precondition (and (next_to ?char ?obj2) (next_to ?char ?obj1) (grabbable ?obj1) (movable ?obj1) (movable ?obj2)))	<b>LLM Output</b> :action lie :precondition (and (next_to ?char ?obj)) (sittable ?obj))	<b>LLM Output</b> :action walk_int :effect (and (inside ?char ?room)) (not (exists (?cur_room - object) (inside ?char ?cur_room))))
Additional Predicates			
<b>Object Property Hallucination</b> <b>Ground Truth</b> :action put_on_character :precondition (or (holds_lh ?char ?obj) (holds_rm ?char ?obj)))	<b>Additional Object State</b> <b>Ground Truth</b> :action clean_stained_brush :parameters (?scrub_brush ?obj) :agent :effects (not (stained ?obj))	<b>Additional Object Orientation</b> <b>Ground Truth</b> :action walk_towards :precondition (and (not (sitting ?char)) (not (lying ?char)))	<b>Additional Spatial Relation</b> <b>Ground Truth</b> :action turn_to :preconditions :effects (facing ?char ?obj)
<b>LLM Output</b> :action put_on_character :precondition (and (holds_rm ?char ?obj) (wearable ?obj)))	<b>LLM Output</b> :action clean_stained_brush :parameters (?scrub_brush ?obj) :agent :effects (and (not (stained ?obj)) (not (stained ?scrub_brush)))	<b>LLM Output</b> :action walk_towards :precondition (and (obj_next_to ?char ?obj) (not (facing ?char ?obj)))	<b>LLM Output</b> :action turn_to :preconditions (next_to ?char ?obj) :effects (facing ?char ?obj)

Figure 29: Transition modeling error examples

1877 that incorporate a broader range of spatial and non-spatial interactions and testing these models in more complex,  
1878 real-world settings. Additionally, exploring hybrid models that combine various strengths of current LLMs could  
1879 yield more capable and flexible systems for practical applications in embodied AI environments.

#### 1880 H.4.3 Error categorization

1881 In assessing the logical form accuracy of predicted PDDL action bodies by the LLM, we categorize errors into  
1882 two main types: missing predicates, and additional predicates as illustrated at figure 29. Each category reflects  
1883 specific discrepancies between the LLM outputs and the ground truth, affecting the precision, recall, and F1  
1884 scores of the model’s predictions.

1885 **Missing Predicates:** These errors occur when essential predicates are omitted in the LLM’s output, leading to  
1886 incomplete or incorrect action specifications. For instance, in the missing object affordance error, the LLM omits  
1887 the ‘(lieable ?obj)’ predicate required for the action LIE, replacing it with incorrect predicates like ‘(sittable  
1888 ?obj)’. Apparently, to lie on some object, the object should be lieable instead of sittable. Such omissions can  
1889 drastically affect the feasibility and correctness of the generated plan, as they fail to capture the necessary  
1890 conditions for action execution.

1891 **Additional Predicates:** This error type is characterized by the inclusion of unnecessary or incorrect predicates  
1892 in the LLM output, which are not present in the ground truth. An example is seen in the additional object state,  
1893 where the LLM predicts ‘(not (stained ?scrub\_brush))’ under effects for the action CLEAN\_STAINED\_BRUSH.  
1894 To clean something with brush, the clean object should be the target object instead of the brush. These additional  
1895 predicates can lead to over-constrained planning scenarios, potentially preventing the execution of valid actions.

1896 Another type of additional predicates is property hallucination. Property hallucination errors arise when the  
1897 LLM inaccurately attributes properties to objects that are not supported by the ground truth. In the depicted  
1898 error, the LLM predicts that an object is ‘(wearable ?obj)’, while ‘wearable’ is not in the vocabulary dictionary.  
1899 This type of error not only introduces factual inaccuracies but also misleads the planning process by enforcing  
1900 constraints that can never be true.

Transition Models all from the Ground Truth - EXECUTION SUCCESS
(:action plug_in :parameters (?char - character ?obj - object) :precondition (or (and (next_to ?char ?obj) (has_plug ?obj) (plugged_out ?obj)) (and (next_to ?char ?obj) (has_switch ?obj) (plugged_out ?obj))) :effect (and (plugged_in ?obj) (not (plugged_out ?obj))) )  (:action walk_towards :parameters (?char - character ?obj - object) :precondition (and (not (sitting ?char)) (not (lying ?char)))

1901

```

(next_to ?char ?obj) (forall (?far_obj - object)
  (when (not (obj_next_to ?far_obj ?obj))
    (not (next_to ?char ?far_obj)))) (forall (?close_obj - object)
  (when (obj_next_to ?close_obj ?obj)
    (next_to ?char ?close_obj))))
:effect (and (next_to ?char ?obj))
)

(:action switch_on
:parameters (?char - character ?obj - object)
:precondition (and (has_switch ?obj) (off ?obj)
  (plugged_in ?obj) (next_to ?char ?obj))
:effect (and (on ?obj) (not (off ?obj)))
)

```

1902

Figure 30: If we use a high-level planner (PDDL planner) over Ground Truth transition models of *plug\_in*, *walk\_towards* and *switch\_on*, the output action sequence can be executed successfully.

1903

1904

#### Transition Models all from LLMs - EXECUTION SUCCESS

```

(:action plug_in
:parameters (?char - character ?obj - object)
:precondition (and (has_plug ?obj) (plugged_out ?obj)
  (next_to ?char ?obj))
:effect (and (plugged_in ?obj) (not (plugged_out ?obj)))
)

(:action walk_towards
:parameters (?char - character ?obj - object)
:precondition (and (inside ?char ?room) (inside ?obj ?room))
:effect (next_to ?char ?obj)
)

(:action switch_on
:parameters (?char - character ?obj - object)
:precondition (and (has_switch ?obj) (off ?obj)
  (next_to ?char ?obj))
:effect (and (on ?obj) (not (off ?obj)))
)

```

1905

Figure 31: If we use a high-level planner (PDDL planner) over GPT-4o predicted transition models of *plug\_in*, *walk\_towards* and *switch\_on*, the output action sequence can be executed successfully.

1906

1907

#### *plug\_in* from LLMs and *walk\_towards* and *switch\_on* from Ground Truth - EXECUTION FAILURE

```

(:action plug_in
:parameters (?char - character ?obj - object)
:precondition (and (has_plug ?obj) (plugged_out ?obj)
  (next_to ?char ?obj))
:effect (and (plugged_in ?obj) (not (plugged_out ?obj)))
)

(:action walk_towards
:parameters (?char - character ?obj - object)
:precondition (and (not (sitting ?char)) (not (lying ?char))
  (next_to ?char ?obj)) (forall (?far_obj - object)
  (when (not (obj_next_to ?far_obj ?obj))
    (not (next_to ?char ?far_obj)))) (forall (?close_obj - object)
  (when (obj_next_to ?close_obj ?obj)
    (next_to ?char ?close_obj))))
:effect (and (next_to ?char ?obj))
)

(:action switch_on
:parameters (?char - character ?obj - object)
:precondition (and (has_switch ?obj) (off ?obj)
  (plugged_in ?obj) (next_to ?char ?obj))
:effect (and (on ?obj) (not (off ?obj)))
)

```

1908

Figure 32: If we use a high-level planner (PDDL planner) over the mixture of transition models, i.e., GPT-4o predicted transition models of *plug\_in* with ground truth *walk\_towards* and *switch\_on*, the output action sequence cannot be executed successfully.

### Consistency of predicted action space

We observe the consistency of LLM predicted action space, that is, the complexity of predicted actions are consistent. Such consistency facilitates the PDDL planner to find possible solutions from initial state to goal state. It is worth noticing that interleaving LLM predicted actions and ground truth actions has lower planner success rate than using only ground truth actions or only LLM predicted actions. Figure 30 demonstrates the success case of PDDL planning using ground truth actions, and figure 31 shows that the actions predicted by GPT-4o also pass the PDDL planner test. However, if we mix the two action spaces, using *plug\_in* from GPT-4o prediction and *walk\_toward* and *switch\_on* from ground truth, the PDDL planner cannot find a feasible solution for this task. This is because *plug\_in* from GPT-4o prediction diverges with *switch\_on* from ground truth. In ground truth action space, to *switch\_on* an object, it needs to be *plugged\_in* and *has\_switch*. Ground truth *plug\_in* can handle the cases when objects either *has\_plug* or *has\_switch*. In the GPT-4o predicted action space, although *plug\_in* cannot handle the cases when objects *has\_switch*, the LLM gets away with it by not specifying *plug\_in* in the preconditions of *switch\_on*. Although the definition is not comprehensive, it can pass the PDDL planner test. However, when mixing the two action spaces together, *switch\_on* requires precondition *plug\_in*. But *plug\_in* cannot handle the cases when objects only *has\_switch* but not *has\_plug*. Thus, the predicted action space maintains consistency, and it is crucial to provide LLM with the context action definitions when we ask LLM to predict a single action.

### Takeaways of Transition Modeling

- (1) **Specialized Performance:** Models such as Claude-3 Opus and Gemini 1.5 Pro excel in specific categories like object states and object orientation, respectively, suggesting that targeted training or specialized architectures enhance LLM capabilities in understanding different types of tasks in transition modeling.
- (2) **Difficulty with Non-Spatial Relations:** Across various models, non-spatial relations consistently pose a challenge, highlighting a gap in the ability of LLMs to grasp complex relational dynamics that are crucial for realistic scenario modeling.
- (3) **Consistency in Action Space:** The effectiveness of planning relies heavily on the consistency of the predicted action space by LLMs; discrepancies between mixed predicted and ground truth actions lead to reduced planner success, emphasizing the need for coherent action definitions in LLM outputs.

## H.5 Correlation with Action Length and Goal Complexity

In this section, we analyze the factors that influence the goal success rates of GPT-4o’s performance in action sequencing for BEHAVIOR . Generally, the goal success rate is influenced by task complexity, as shown in Figure 33. The number of goals within a task, the number of state goals, and the number of relation goals all adversely affect the success rates. The length of the ground-truth action sequence follows the same trend, although there are some fluctuations.

Success rates for different actions and predicates are provided in Figure 34. The predicate with the highest success rate is *open*, which corresponds to the action *OPEN* that has a 1.00 execution success rate. This high success rate is because GPT-4o can successfully reason about the preconditions for this action, as shown in Table 24.

The complexity of action preconditions also adversely affects the success rate. Actions with lower execution success rates often require tool usage, which is a challenge for LLMs. For example, *Soak* requires placing the object in a toggled-on sink, *Slice* requires holding a knife, and *Clean* requires a cleaning tool like a rag or a brush, and if cleaning a stain, the cleaning tool needs to be soaked.

Interestingly, the success rate of the same actions performed with the right hand is slightly higher than those performed with the left hand. The action with the lowest success rate is *SLICE*. There are only two cases in BEHAVIOR that require this action, and GPT-4o failed both.

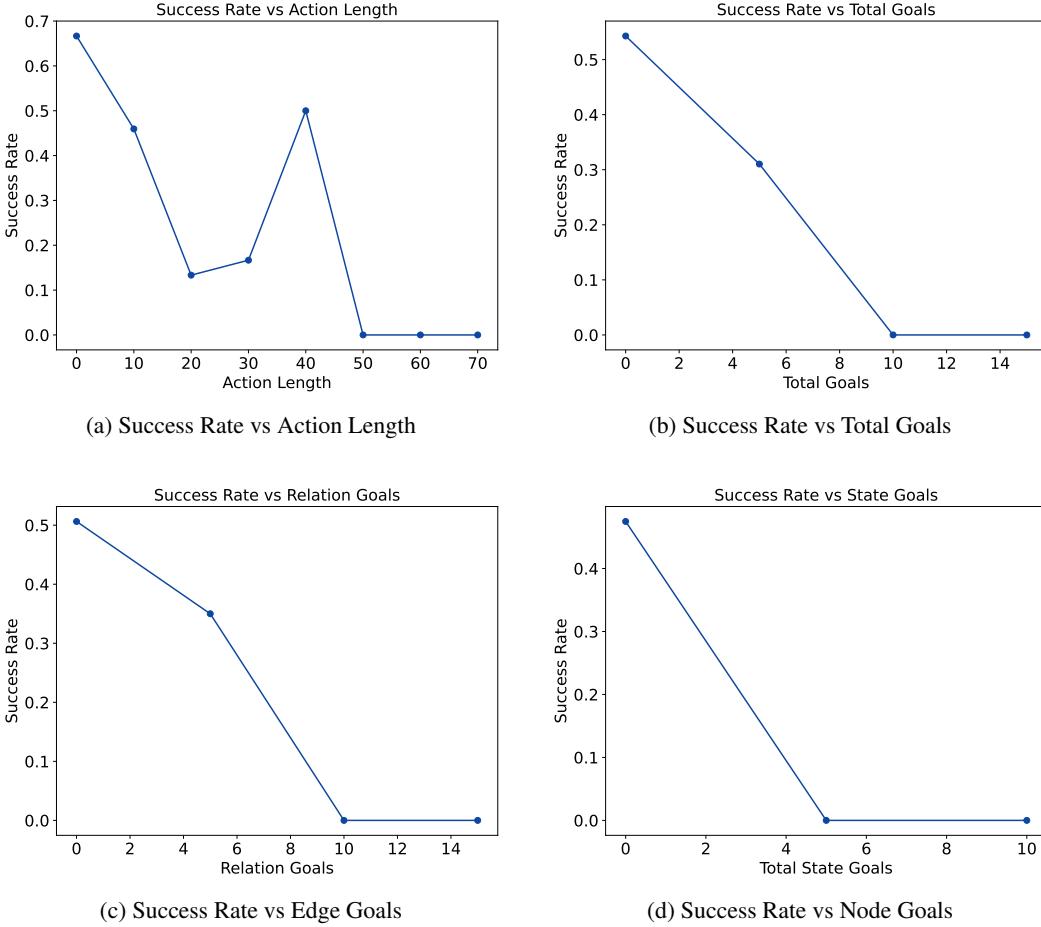


Figure 33: BEHAVIOR Success Rate as a Function of Different Cofactors

#### Takeaways of Error Factor Analysis

1. Task complexity adversely affects goal success rates, including the number of total goals, state goals, relation goals, and ground-truth action sequence length.
2. Complexity of action preconditions adversely affects success rates, especially for actions requiring tool usage, e.g., *Soak*, *Slice*, and *Clean*.
3. Actions performed with the right hand have a slightly higher success rate than those with the left hand.

1946

## 1947 I Sensitivity Analysis

### 1948 I.1 Motivation and Problem Formulation

1949 In the current setting of transition modeling task, we use the entire space of LLM predicted actions combined  
 1950 with PDDL planner to check whether there exists a feasible solution fulfilling the goals. However, it still lacks  
 1951 finer-grid metrics to show which action fails if the PDDL planner fails to find a solution. Therefore, we conduct  
 1952 sensitivity analysis to examine how sensitive task success rates are to specific LLM predicted actions. In the  
 1953 ground truth space, if after replacing a specific action by the LLM predicted one, the PDDL planner fails to find  
 1954 a solution, we call the task is *sensitive* to this action, and vice versa. This sensitivity analysis provides insights  
 1955 into the per action impact of LLM predictions on the success rate of specific actions within each task category,  
 1956 highlighting areas where the LLM performs well and where it may require further improvement. Sensitivity  
 1957 analysis also highlights the actions that are crucial in different task categories, providing insights for fine-tuning  
 1958 and downstream tasks.

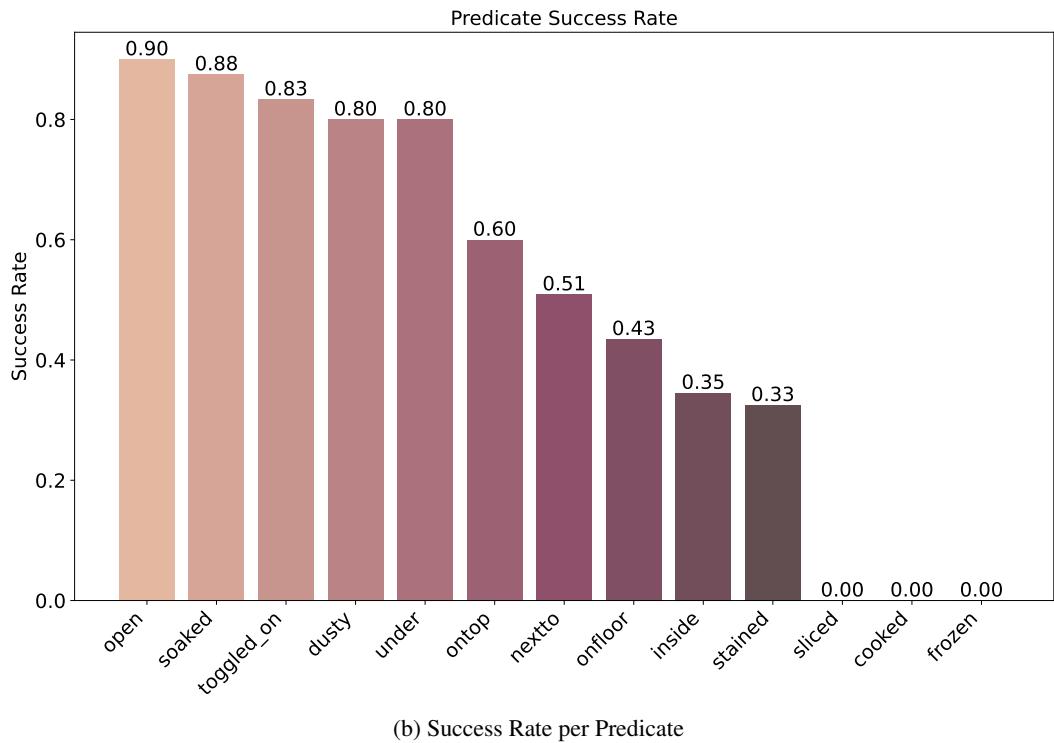
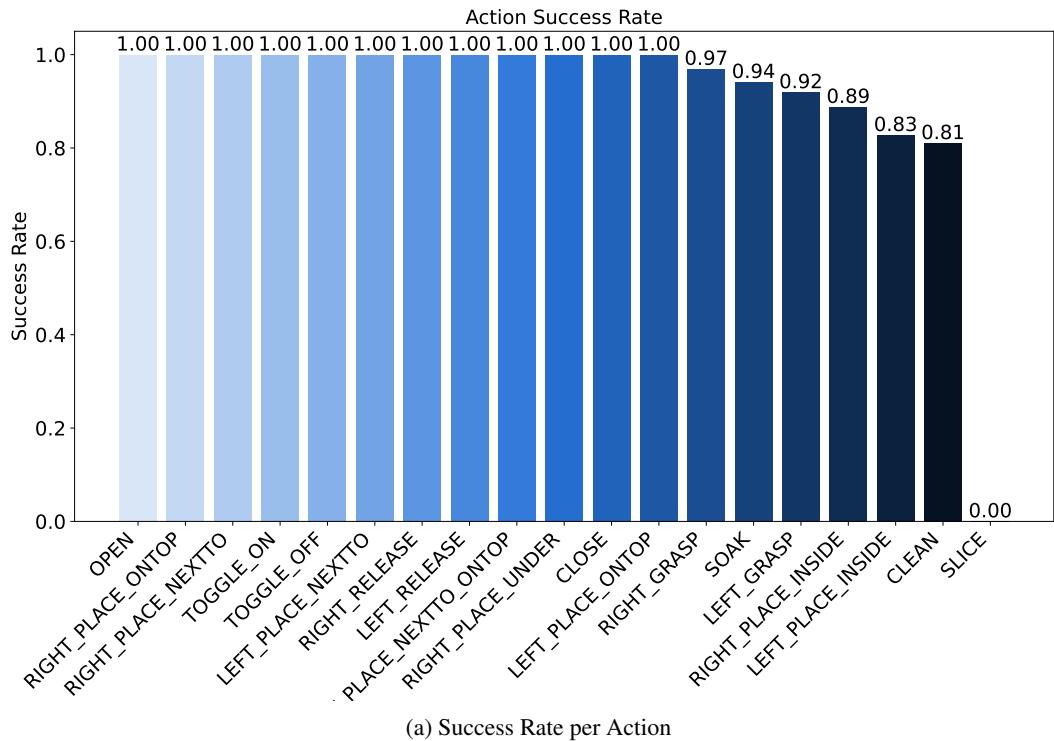


Figure 34: BEHAVIOR Success Rate Analysis for Actions and Predicates

1959 **I.2 Implementation Details**

1960 We start the sensitivity analysis from setting the action space as the ground truth space. Given the task, for  
 1961 each relevant action, we replace one ground truth action by the LLM predicted counterpart, and solve the task  
 1962 by PDDL planner. Each task is categorized according to Appendix M.6 and we report the overall sensitivity  
 1963 analysis results and the results by task categories.

1964 **I.3 Result Analysis**

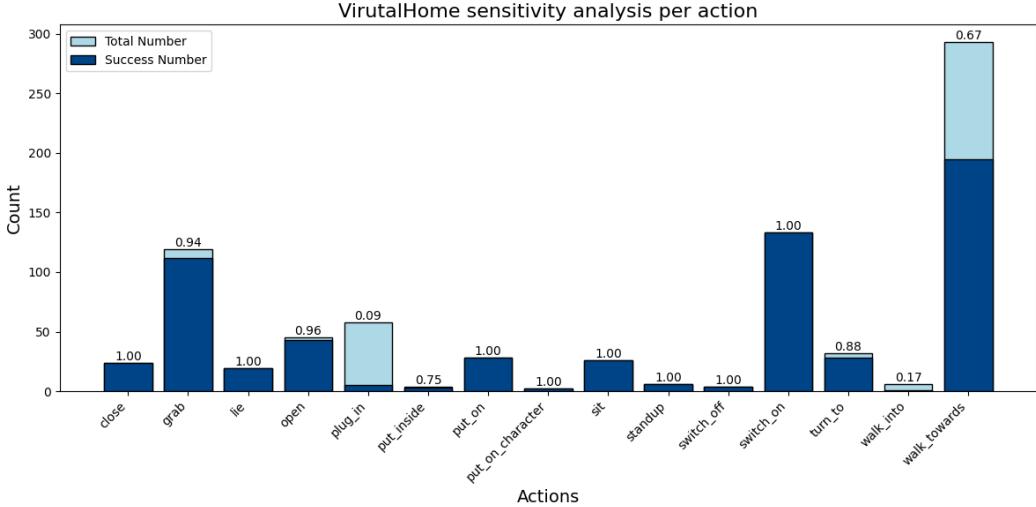


Figure 35: VirtualHome sensitivity analysis per action

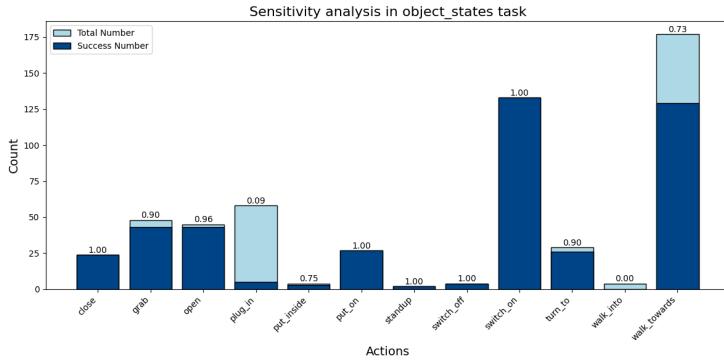


Figure 36: Sensitivity analysis for object states tasks in VirtualHome

1965 The sensitivity analysis conducted on a large language model's predictions for various tasks provides a detailed  
 1966 view of the model's robustness and areas necessitating improvement. This analysis evaluates the impact of  
 1967 replacing specific actions predicted by the model while maintaining other actions as ground truth across five  
 1968 different task categories: non-spatial relations, object affordance, object orientation, object states, and spatial  
 1969 relations.

1970 In the *non-spatial relations* task, the model exhibited a nearly perfect success rate for most actions except for  
 1971 "grab" and "walk\_towards," with success rates of 0.97 and 0.65, respectively. The high success rates indicate a  
 1972 strong model performance in scenarios where spatial relationships are not a factor, although the lower rate for  
 1973 "walk\_towards" suggests difficulties in predicting movement-related actions when spatial context is absent.

1974 For the *object affordance* task, there were notable variances, with "plug\_in" and "put\_inside" having significantly  
 1975 lower success rates of 0.09 and 0.75. This disparity suggests challenges in the model's understanding of actions  
 1976 involving complex interactions or manipulations, whereas actions directly associated with straightforward  
 1977 affordances like "close" and "grab" achieved perfect success.

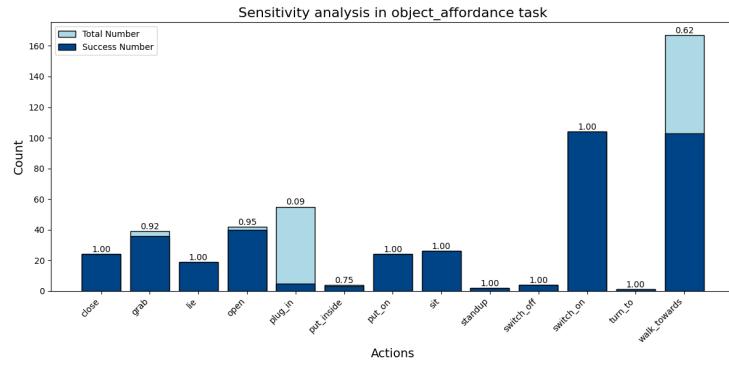


Figure 37: Sensitivity analysis for object affordance tasks in VirtualHome

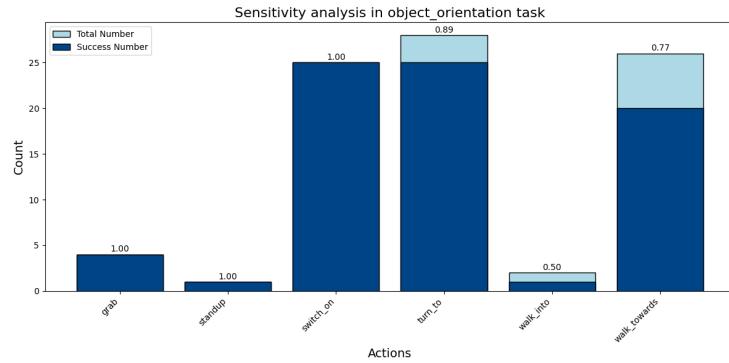


Figure 38: Sensitivity analysis for object orientation tasks in VirtualHome

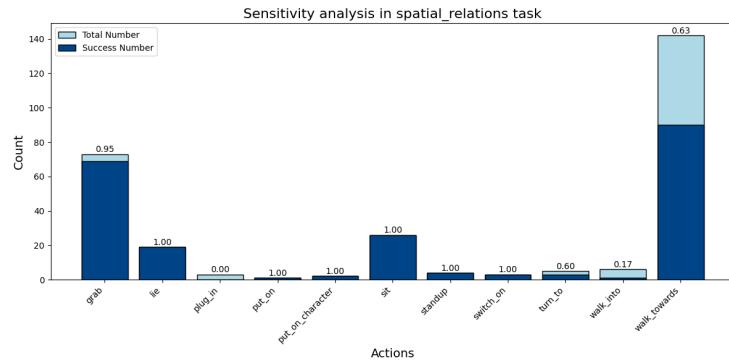


Figure 39: Sensitivity analysis for spatial relations tasks in VirtualHome

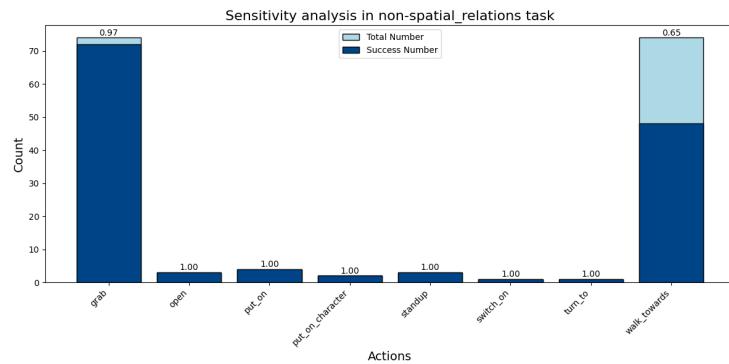


Figure 40: Sensitivity analysis for non-spatial relations tasks in VirtualHome

- 1978 The *object orientation* task showed high success rates for most actions, but "turn\_to" and "walk\_into" had lower  
 1979 success rates of 0.89 and 0.77. This indicates that the model is proficient with simple orientation changes yet  
 1980 slightly struggles with actions requiring precise orientation or movement towards specific objects.
- 1981 In the *object states* task, "plug\_in" displayed extremely low success rates of 0.09. The action "plug\_in" likely  
 1982 presents conceptual challenges or ambiguities in interpretation, pointing to a critical need for enhancing the  
 1983 model's training on the diverse contexts and functionalities of objects.
- 1984 The *spatial relations* task results highlight significant challenges, as seen with "walk\_towards" and "walk\_into"  
 1985 having success rates of 0.63 and 0.17. These outcomes underscore difficulties in handling spatially complex  
 1986 predictions and interactions within a spatial context.
- 1987 From Figure 35, generally, we notice "plug\_in" and "walk\_towards" as the most challenging actions for LLM to  
 1988 predict. The preconditions of "plug\_in" in VirtualHome contains two cases: either the device "has\_plug", or the  
 1989 device "has\_switch". While most LLMs are able to predict "has\_plug", they can hardly catch the other case,  
 1990 which turns out to be quite common in the tasks. The difficulty of "walk\_towards" lies in its complex constraints  
 1991 on spatial relations. Many LLMs either predict additional constraints or miss some necessary constraints, making  
 1992 the success rate of "walk\_towards" low.

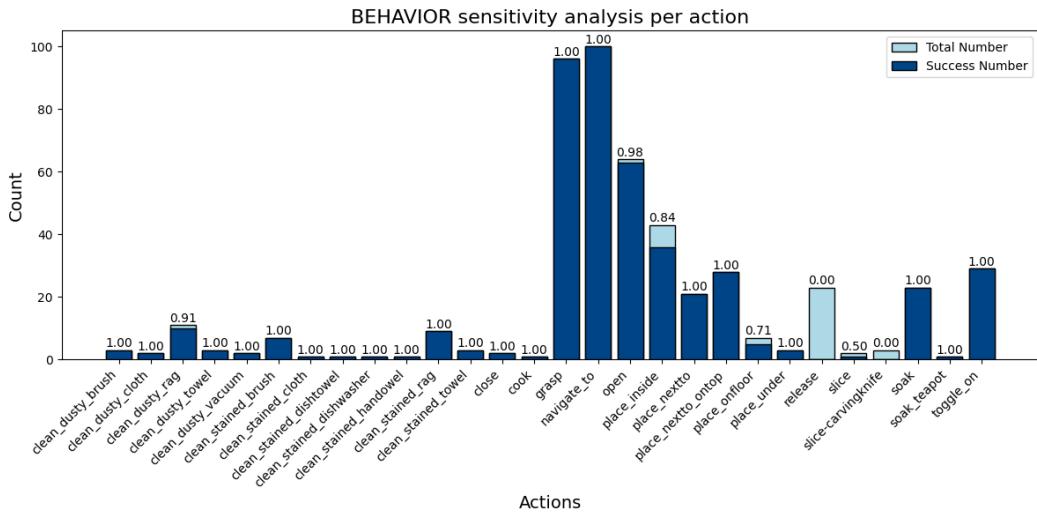


Figure 41: BEHAVIOR sensitivity analysis per action

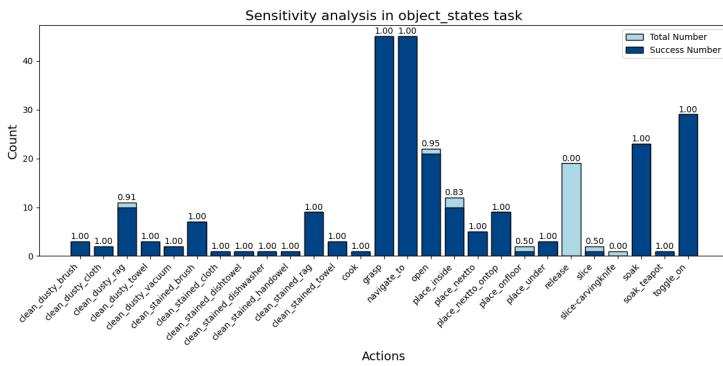


Figure 42: Sensitivity analysis for object states tasks in BEHAVIOR

- 1993 The analysis conducted on the BEHAVIOR dataset for the tasks of non-spatial relations, object states, and  
 1994 spatial relations offers insights into the predictive capabilities of the LLM and delineates areas that necessitate  
 1995 improvements.
- 1996 In the *non-spatial relations* task, the analysis revealed that the model performs exceptionally well for  
 1997 most actions, achieving perfect success rates. However, deviations are observed in the actions "release"  
 1998 and "slice\_carvingknife" with success rates of 0.00 and 0.00, respectively. The low success rate for

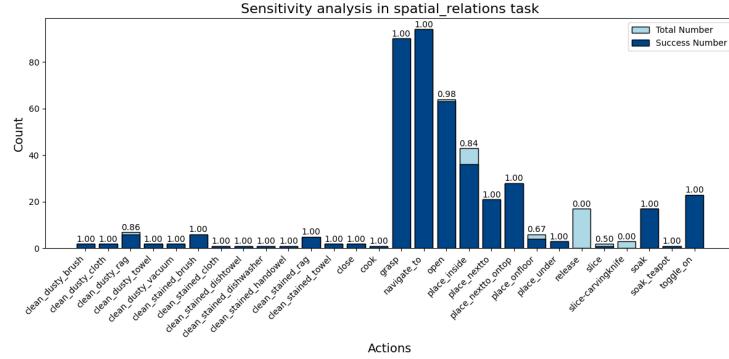


Figure 43: Sensitivity analysis for spatial relations tasks in BEHAVIOR

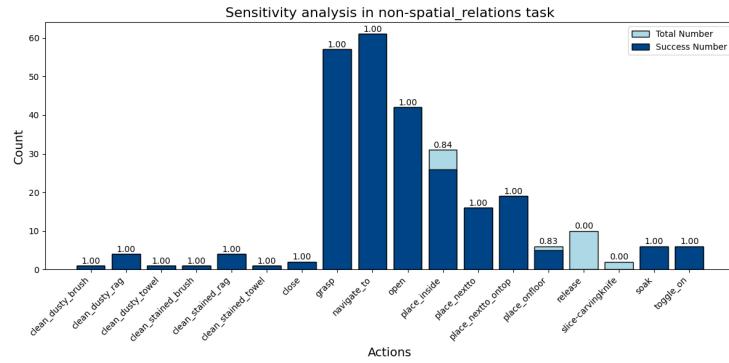


Figure 44: Sensitivity analysis for non-spatial relations tasks in BEHAVIOR

- 1999 "slice\_carvingknife" suggests that the model struggles significantly with actions that involve complex interactions, including tool use and precise movements, potentially due to inadequate representation in the training data.
- 2000 For the *object states* task, similar trends are evident where actions such as "grasp" and "release" exhibit lower success rates of 0.95 and 0.83, respectively, and "slice" again recorded a success rate of 0.00. These findings indicate that tasks requiring detailed manipulation of object states are challenging for the model, pointing to a gap in training on nuanced and precise operations.
- 2001 The *spatial relations* task also highlights robust performances in most actions, yet "grasp" and "slice" show lower success rates of 0.84 and 0.00. The consistent underperformance on "slice" across different tasks underscores a significant deficiency in the model's ability to handle precise tool-based actions, while the challenges with "grasp" suggest difficulties in scenarios that demand fine motor skills or specific spatial awareness.
- 2002 Across all categories in the BEHAVIOR dataset, as shown in figure 41, the model's proficiency in straightforward actions contrasts with its struggles in more complex or nuanced actions. For example, "place\_inside" shows rather low success rate in many tasks, probably because it requires understanding of both spatial relations like "inside", and non-spatial relations like "holding". Agent also needs to be careful whether the robot currently is "handsfull". All of these decisions are challenging for current LLMs. Generally, figure 41 underscores a limitation in the model's current training, which may not sufficiently encompass the complexity of real-world interactions or the specificity required for precise task execution.
- 2003 Overall, the model performs commendably on tasks involving basic manipulations but shows limitations in complex spatial judgments or detailed object interactions. To address these deficiencies, it is recommended to enhance the model's training with more diverse scenarios focusing on underperforming actions, fine-tune using specialized datasets, and conduct a thorough error analysis to precisely identify and rectify the sources of errors.
- 2004 This strategy will not only improve the model's performance but also expand its applicability across varied tasks, reinforcing its utility in complex real-world applications.

## Takeaways of Sensitivity Analysis

- (1) **Action-Specific Sensitivities:** The analysis reveals significant variations in the success rates for specific actions across tasks, with actions like "plug\_in" and "walk\_towards" consistently showing low success due to complex preconditions and spatial requirements. This underscores the need for LLMs to better understand and generate actions with multi-faceted constraints.
- (2) **Challenges with Complex Interactions:** Complex interactions involving detailed object manipulation, such as "slice\_carvingknife" and "place\_inside," present notable challenges, highlighting deficiencies in the model's training on precise movements and tool use. Enhancing data representation for such scenarios in training sets is crucial for improving performance.
- (3) **Robustness in Basic Tasks vs. Nuanced Deficiencies:** While the model performs well in straightforward tasks, it struggles with nuanced actions that require fine motor skills or advanced spatial awareness. This contrast points to the necessity of incorporating more complex and varied interaction scenarios in training to refine the model's capabilities in handling real-world complexities.
- (4) **Training and Fine-Tuning Needs:** The consistent issues across different types of tasks suggest that the current training regimens may not fully capture the diversity of real-world interactions, especially in spatial and object-oriented tasks. A targeted approach to fine-tuning the model with specialized datasets that focus on underperforming actions could lead to significant improvements in model robustness and applicability in practical settings.

2023

## 2024 J Pipeline-Based vs Modularized

### 2025 J.1 Motivation and Problem Formulation

2026 **Motivation** In previous sections, we examined four modularized ability modules in EMBODIED AGENT INTERFACE. Despite covering all four abilities, it remains important to understand how LLMs address complete tasks  
2027 when provided with natural goals instead of symbolic ones. This is crucial because, in most settings, only natural  
2028 language goal descriptions are available, and symbolic goals are often missing. Natural language descriptions  
2029 can be ambiguous. For instance, the natural language instruction *cook food* can be as simple as *cook the*  
2030 *chicken then eat it*, whereas its symbolic goal should include detailed information such as *cooked(chicken)*  
2031 and *ontop(chicken, plate)* and *ontop(plate, table)*. One advantage of our interface is  
2032 the support for flexible composition of primitive modules. This allows us to first have LLMs parse natural  
2033 language goals into symbolic goals using the goal interpretation ( $\mathcal{G}$ ) module before conducting further analysis.  
2034 Therefore, we implement pipeline-based methods to further investigate such ability of LLMs.  
2035

2036 **Problem Formulation** Given an initial state  $s_0$  and a natural language goal  $g_{nl}$ , our objective is to have LLMs  
2037 generate a feasible action sequence  $\bar{a}$  to achieve  $g_{nl}$  and to decompose  $g_{nl}$  into a subgoal trajectory  $\bar{\phi}$ .

### 2038 J.2 Implementation Details

2039 We evaluate LLM abilities through pipeline-based methods in BEHAVIOR enviroment. Using the three  
2040 methods—Goal Interpretation ( $\mathcal{G}$ ), Action Sequencing ( $\mathcal{Q}$ ), and Subgoal Decomposition ( $\Phi$ ) provided in  
2041 EMBODIED AGENT INTERFACE, we can establish two pipelines to obtain  $\bar{a}$  and  $\bar{\phi}$  respectively. Specifically,  
2042 by leveraging the rule  $\mathcal{G} : \langle s_0, g_{nl} \rangle \rightarrow g$ , we derive the symbolic goal  $g$  corresponding to  $g_{nl}$ . Then, we (1)  
2043 apply the rule  $\mathcal{Q} : \langle s_0, g \rangle \rightarrow \bar{a}$  to get the action sequence  $\bar{a}$ , and (2) apply the rule  $\Phi : \langle s_0, g \rangle \rightarrow \bar{\phi}$  to derive the  
2044 subgoal trajectory  $\bar{\phi}$ .

### 2045 J.3 Result Analysis

2046 Table 18 presents a comparison of modularized methods and pipeline-based methods for pipelines  $\mathcal{G} + \mathcal{Q}$  and  
2047  $\mathcal{G} + \Phi$ . Our observations indicate the following: (1) The trajectory executable rates are similar between both  
2048 methods. (2) Pipeline-based methods suffer from error accumulation due to the composition of two modules.

2049 These findings are consistent across both pipelines  $\mathcal{G} + \mathcal{Q}$  and  $\mathcal{G} + \Phi$ . We attribute our observations to two  
2050 main reasons: First, achieving an executable trajectory requires an understanding of semantics, which is largely  
2051 related to the provided environment and vocabulary. Since these factors remain unchanged between modularized  
2052 and pipeline-based methods, Language Models (LLMs) generate comparable executable trajectories for both.  
2053 Second, pipeline-based methods tend to perform worse in terms of goal success rate. Errors made by the goal  
2054 interpretation model mislead the LLMs in downstream modules, causing them to generate less accurate action  
2055 and subgoal sequences aimed at the goal. Nevertheless, despite these challenges, LLMs are still capable of  
2056 generating feasible and executable action sequences.

2057 Furthermore, Table 18 demonstrates that while most state-of-the-art (SOTA) LLMs tend to avoid grammar errors.  
2058 However, this finding does not hold for less advanced models, including Gemini 1.0 Pro, GPT-3.5-turbo, and  
2059 some open-sourced LLMs. We attribute this to the fact that most SOTA LLMs are fine-tuned to follow human

2060 instructions more closely, whereas less advanced models either have smaller sizes or undergo less rigorous  
 2061 tuning strategies. Regarding runtime errors, all LLMs, regardless of their advancement, tend to miss necessary  
 2062 steps in their generation process. This observation also aligns with our findings discussed in Section H.2.

### Takeaways of Pipeline-Based vs Modularized

- (1) **Similar Trajectory Execution Rates:** Both modularized and pipeline-based methods have similar trajectory executable rates for pipelines  $\mathcal{G} + \mathcal{Q}$  and  $\mathcal{G} + \Phi$ .
- (2) **Error Accumulation in Pipelines:** Pipeline-based methods suffer from error accumulation, particularly due to the composition of two modules, which impacts their performance negatively.
- (3) **Grammar Error Trends:** State-of-the-art (SOTA) LLMs generally avoid grammar errors for both pipeline-based and modularized methods, unlike less advanced models such as Gemini 1.0 Pro, GPT-3.5-turbo, and some open-sourced LLMs.
- (4) **Runtime Errors Persist:** All LLMs, regardless of their level of advancement, are prone to runtime errors, missing necessary steps in their generation process.

2063

Table 18: Pipeline-based evaluation results for (1) $\mathcal{G} + \mathcal{Q}$  and (2)  $\mathcal{G} + \Phi$  in BEHAVIOR .  $\mathcal{G}$ : Goal Interpretation.  $\mathcal{Q}$ : Action Sequencing.  $\Phi$ : Subgoal Decomposition. In this table, M means ‘modularized’, whereas P means ‘pipeline-based’.

Model	Goal Evaluation				Trajectory Evaluation														
	Goal SR		Execution SR		Grammar Error ( $\downarrow$ )			Runtime Error ( $\downarrow$ )											
	M	P	M	P	M	P	M	P	M	P	M	P	M	P	M	P			
<i>Goal Interpretation + Action Sequencing</i>																			
Claude-3 Haiku	26.0	21.0	32.0	29.0	0.0	0.0	6.0	6.0	0.0	0.0	7.0	6.0	54.0	52.0	1.0	7.0	1.0	17.0	
Claude-3 Sonnet	44.0	41.0	57.0	53.0	0.0	0.0	1.0	3.0	0.0	0.0	11.0	14.0	19.0	21.0	11.0	9.0	2.0	12.0	
Claude-3 Opus	51.0	46.0	59.0	54.0	0.0	1.0	0.0	1.0	0.0	0.0	3.0	6.0	35.0	35.0	3.0	3.0	2.0	4.0	
Gemini 1.0 Pro	27.0	26.0	32.0	35.0	7.0	5.0	3.0	3.0	6.0	6.0	13.0	14.0	35.0	38.0	4.0	2.0	4.0	11.0	
Gemini 1.5 Flash	40.0	35.0	52.0	49.0	0.0	0.0	0.0	2.0	0.0	0.0	5.0	10.0	42.0	41.0	1.0	0.0	2.0	7.0	
Gemini 1.5 Pro	42.0	37.0	54.0	55.0	0.0	1.0	0.0	1.0	0.0	0.0	6.0	7.0	39.0	35.0	1.0	1.0	2.0	0.0	
GPT-3.5-turbo	16.0	14.0	20.0	32.0	4.0	1.0	7.0	3.0	23.0	15.0	1.0	5.0	36.0	39.0	8.0	6.0	1.0	3.0	
GPT-4-turbo	38.0	32.0	45.0	47.0	0.0	1.0	0.0	1.0	0.0	0.0	7.0	9.0	47.0	41.0	1.0	1.0	0.0	0.0	
GPT-4o	47.0	42.0	53.0	55.0	0.0	0.0	1.0	3.0	0.0	0.0	9.0	6.0	36.0	35.0	1.0	1.0	0.0	4.0	
Cohere Command R	16.0	5.0	19.0	9.0	5.0	3.0	13.0	38.0	0.0	1.0	8.0	8.0	43.0	31.0	12.0	12.0	4.0	8.0	
Cohere Command R+	27.0	15.0	35.0	29.0	0.0	0.0	1.0	8.0	15.0	14.0	10.0	10.0	30.0	39.0	31.0	0.0	2.0	15.0	
Mistral Large	33.0	31.0	50.0	38.0	0.0	0.0	0.0	3.0	0.0	0.0	8.0	14.0	35.0	37.0	6.0	8.0	7.0	5.0	
Mistral 8x22B MoE	30.0	26.0	40.0	36.0	3.0	3.0	6.0	13.0	0.0	0.0	10.0	14.0	32.0	21.0	9.0	13.0	2.0	15.0	
Llama3 8B	10.0	0.0	16.0	5.0	0.0	2.0	15.0	25.0	9.0	6.0	11.0	44.0	34.0	9.0	17.0	5.0	14.0		
Llama3 70B	34.0	26.0	42.0	40.0	0.0	1.0	2.0	3.0	0.0	0.0	15.0	18.0	38.0	35.0	3.0	5.0	6.0	9.0	
<i>Goal Interpretation + Subgoal Decomposition</i>																			
Claude-3 Haiku	29.0	21.0	35.0	40.0	0.0	0.0	1.0	5.0	0.0	0.0	2.0	2.0	59.0	46.0	3.0	7.0	3.0	16.0	
Claude-3 Sonnet	38.0	31.0	43.0	45.0	0.0	0.0	2.0	3.0	0.0	0.0	3.0	2.0	51.0	47.0	1.0	3.0	3.0	18.0	
Claude-3 Opus	39.0	35.0	47.0	45.0	0.0	0.0	3.0	8.0	0.0	0.0	5.0	4.0	45.0	42.0	0.0	1.0	5.0	7.0	
Gemini 1.0 Pro	23.0	14.0	33.0	30.0	2.0	0.0	4.0	10.0	0.0	1.0	3.0	3.0	1.0	51.0	45.0	7.0	13.0	3.0	17.0
Gemini 1.5 Flash	34.0	32.0	42.0	44.0	2.0	1.0	1.0	3.0	0.0	0.0	2.0	2.0	53.0	48.0	0.0	2.0	3.0	7.0	
Gemini 1.5 Pro	31.0	26.0	37.0	38.0	0.0	1.0	1.0	3.0	0.0	0.0	3.0	2.0	59.0	56.0	0.0	0.0	2.0	1.0	
GPT-3.5-turbo	24.0	14.0	36.0	27.0	2.0	0.0	3.0	12.0	0.0	22.0	3.0	1.0	52.0	32.0	4.0	6.0	3.0	5.0	
GPT-4-turbo	37.0	37.0	47.0	49.0	0.0	0.0	3.0	4.0	0.0	0.0	9.0	8.0	40.0	37.0	1.0	2.0	6.0	6.0	
GPT-4o	48.0	38.0	55.0	52.0	0.0	0.0	3.0	4.0	0.0	0.0	5.0	6.0	37.0	35.0	0.0	3.0	5.0	9.0	
Cohere Command R	15.0	8.0	25.0	15.0	21.0	13.0	11.0	32.0	0.0	1.0	1.0	1.0	38.0	32.0	4.0	6.0	4.0	12.0	
Cohere Command R+	24.0	17.0	37.0	31.0	2.0	6.0	4.0	10.0	0.0	2.0	5.0	7.0	51.0	40.0	1.0	4.0	6.0	14.0	
Mistral Large	30.0	22.0	38.0	29.0	1.0	1.0	3.0	12.0	0.0	1.0	4.0	5.0	52.0	50.0	2.0	2.0	1.0	5.0	
Mistral 8x22B MoE	27.0	22.0	33.0	29.0	0.0	0.0	4.0	9.0	0.0	2.0	2.0	2.0	59.0	45.0	2.0	13.0	0.0	17.0	
Llama3 8B	21.0	3.0	29.0	14.0	2.0	7.0	11.0	29.0	0.0	2.0	6.0	3.0	44.0	30.0	8.0	15.0	7.0	7.0	
Llama3 70B	20.0	19.0	30.0	31.0	1.0	1.0	5.0	22.0	1.0	1.0	8.0	7.0	51.0	35.0	4.0	3.0	4.0	7.0	

## K Replanning and Feedback

### K.1 Motivation and Problem Formulation

In our current setting, the agent has only one chance to generate plans or make prediction regardless of previous failure and warning. However, in reality, it is an important ability for agent to learn from feedback and re-plan based on the failure and warning. It demonstrates the agent’s ability to make agile adjustments based on the simulator execution and environment. Practically, replanning and feedback also helps to prevent the agent from making same mistakes over and over again. Therefore, we regard model’s ability to re-plan from feedback necessary and vital.[7][72][73]

### K.2 Implementation Details

To evaluate agent’s ability to re-plan from the feedback, for each unsuccessful task, either failed in the execution (grammar error, runtime error etc.), or failed in the goal satisfaction check (unsatisfied state goals, relation goals,

2075 or action goals), we construct the feedback as the error message and necessary information, and append them  
 2076 to the task prompt. We provide the previous agent's plan as "*At the [retry\_cnt] retry, LLM predict the action*  
 2077 *sequence to be [predicted\_action]*", and detail the reason why it fails. For example, for runtime error missing  
 2078 step, the error message will be *Action [action] is not executable in the action sequence [actions]. It encounters*  
 2079 *error: MISSING STEP. Missing step means that action [action] needs some other necessary action before its*  
 2080 *execution.*". Another error message example for unsatisfied goals is *Action sequence [actions] does not satisfy*  
 2081 *all the goals. Please check the action sequence and try again. Specifically, the following goals are not satisfied:*  
 2082 *Node goals not satisfied: [unsatisfied\_node\_goals] Edge goals not satisfied: [unsatisfied\_edge\_goals] Action*  
 2083 *goals not satisfied: [unsatisfied\_action\_goals].* We allow agent to re-plan for at most 3 times, and each time, we  
 2084 append all previous attempts and error messages to the feedback. We test this setting on the action sequencing  
 2085 task with model GPT-4o.

2086 **K.3 Result Analysis**

Table 19: Replanning evaluation results (%) for *action sequencing*

Model	Goal Evaluation		Trajectory Evaluation					
	Goal SR	Execution SR	Grammar Error (↓)			Runtime Error (↓)		
			Parsing	Hallucination	Action-Arg Num	Wrong Order	Missing Step	Affordance
GPT-4o	65.2	71.8	0.0	1.3	0.7	0.0	25.3	1.0
GPT-4o w/ replanning	77.4	83.3	0.0	1.3	0.0	0.0	14.1	0.3
								0.7

2087 From table 19, we observe that with replanning, the model gains improvement by more than 10% on the success  
 2088 rate and executable success rate. Other metrics also gain improvement except the additional steps error rate.  
 2089 Possible reason could be that agent sometimes tries to iterate over previous attempts and make modification  
 2090 based on those attempts. Such strategy makes agent easily overstate some actions and cause additional step  
 2091 warning.

**Takeaways of Replanning and Feedback**

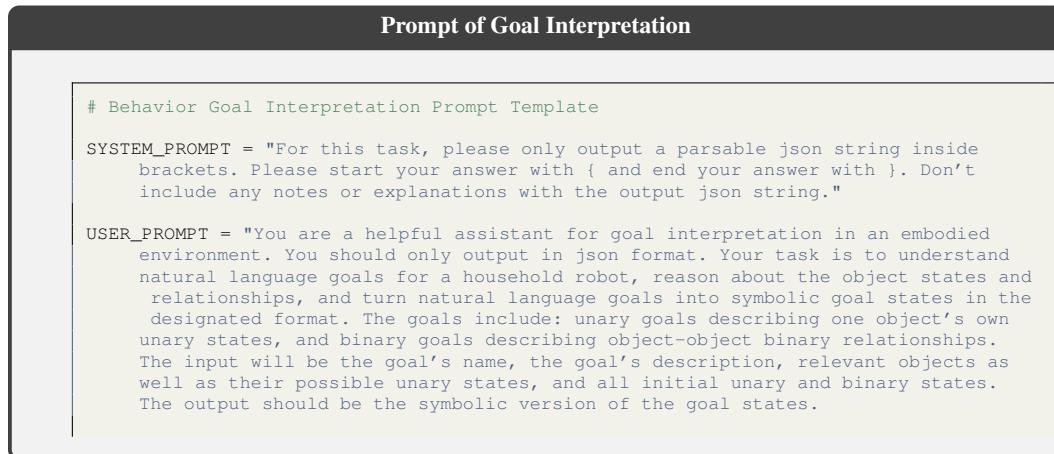
- (1) **Enhanced Adaptability through Replanning:** Incorporating the ability for models to replan based on feedback significantly improves their performance, demonstrating over a 10% increase in success rates. This shows that LLMs can effectively learn from their errors and adjust their action sequences to better align with the required task outcomes, enhancing overall adaptability. Replanning also helps to reduce the repetition of similar errors, showcasing the model's ability to evolve its strategy over successive iterations.
- (2) **Risk of Over-Correction:** While replanning generally leads to improvements in task execution, it can sometimes result in the over-generation of actions, as indicated by an increased rate of additional steps errors. This suggests that while the model is adept at adjusting to feedback, it may also benefit from more precise guidance to avoid introducing unnecessary actions in its revised plans.

2092

2093 **L Prompt and Analysis**

2094 **L.1 Prompt of Goal Interpretation**

2095 We design the prompt template of goal interpretation as below:



2096

2097

2098  
2099

## 2100 L.2 Prompt of Subgoal Decomposition

2101 We design the prompt template of subgoal decomposition as below:

### Prompt of Subgoal Decomposition

```
# Behavior Subgoal Decomposition Prompt Template

SYSTEM_PROMPT = '''# Background Introduction
You are determining the complete state transitions of a household task solving by
a robot. The goal is to list all intermediate states (which is called subgoals
as a whole) to achieve the final goal states from the initial states. The output
consists of a list of boolean expression, which is a combination of the state
predicates. Note that your output boolean expression list is in temporal order,
therefore, it must be consistent and logical. In short, your task is to output
the subgoal plan in the required format.

# Data Vocabulary Introduction
Below we introduce the detailed data vocabulary and their format that you can use
to generate the subgoal plan.
## Available States
The state is represented as a first-order predicate, which is a tuple of a
predicate name and its arguments. Its formal definition looks like this "<
PredicateName>(Params)", where <PredicateName> is the state name and each param
should be ended with an id. An example is "inside(coin.1, jar.1)". Below is a
list of available states and their descriptions.
<available_behavior_states>
## Available Connectives
```

2102

```

The connectives are used to satisfy the complex conditions. They are used to
combine the state predicates.
<available_behavior_connectives>
# Rules You Must Follow
- The initial states are the states that are given at the beginning of the task.
- Your output must be a list of boolean expressions that are in temporal order.
- You must follow the data format and the available states and connectives
defined above.
- The output must be consistent, logical and as detailed as possible. View your
output as a complete state transition from the initial states to the final goal
states.

- Please note that the robot can only hold one object in one hand. Also, the
robot needs to have at least one hand free to perform any action other than put,
place, take, hold.
- Use holds_rh and holds_lh in your plan if necessary. For example, stained(shoe)
cannot directly change to not stained(shoe), but needs intermediate states like
[soaked(rag), holds_rh(rag), not stained(shoe)] or [holds_rh(detergent), not
stained(shoe)].
- Your output follows the temporal order line by line. If you think there is no
temporal order requirement for certain states, you can use connective "and" to
combine them. If you think some states are equivalent, you can use connective "
or" to combine them.
- Please use provided relevant objects well to help you achieve the final goal
states. Note that inside(obj1, agent) is an invalid state, therefore you cannot
output it in your plan.
- Do not output redundant states. A redundant state means a state that is either
not necessary or has been satisfied before without broken.
- Your must strictly follow the json format like this {"output": [<your subgoal
plan>]}, where <your subgoal plan> is a list of boolean expressions presented in
the temporal order.
- Start your output with "{" and end with "}". For each line of the output, DO
NOT INCLUDE IRRELEVANT INFORMATION (like number of line, explanation, etc.).

# Example: Task is bottling_fruit
Below we provide an example for your better understanding.
<in_context_example>''

USER_PROMPT = '''Now, it is time for you to generate the subgoal plan for the
following task.
# Target Task: <task_name>
## Relevant objects in this scene
<relevant_objects>
## Initial States
<initial_states>
## Goal States
<goal_states>
## Output: Based on initial states in this task, achieve final goal states
logically and reasonably. It does not matter which state should be satisfied
first, as long as all goal states can be satisfied at the end. Make sure your
output follows the json format, and do not include irrelevant information, do
not include any explanation. Output concrete states and do not use quantifiers
like "forall" or "exists".'''


```

2103

2104

2105

### 2106 L.3 Prompt of Action Sequencing

2107 We design the prompt template of action sequencing as below:

**Prompt of Action Sequencing**

```

SYSTEM_PROMPT = """For this task, please only output a parsable json string inside
brackets. Please start your answer with [ and end your answer with ]. Don't
include any notes or explanations with the output json string."""

USER_PROMPT=""

Problem:
You are designing instructions for a household robot.
The goal is to guide the robot to modify its environment from an initial state to a
desired final state.

```

2108

The input will be the initial environment state, the target environment state, the objects you can interact with in the environment.

The output should be a list of action commands so that after the robot executes the action commands sequentially, the environment will change from the initial state to the target state.

Data format: After # is the explanation.

Format of the states:  
The environment state is a list starts with a unary predicate or a binary predicate, followed by one or two objects.  
You will be provided with multiple environment states as the initial state and the target state.

For example:

```
['inside', 'strawberry_0', 'fridge_97'] #strawberry_0 is inside fridge_97
[not', 'sliced', 'peach_0'] #peach_0 is not sliced
['ontop', 'jar_1', 'countertop_84'] #jar_1 is on top of countertop_84
```

Format of the action commands:  
Action commands is a dictionary with the following format:

```
{}  
    'action': 'action_name',  
    'object': 'target_obj_name',  
}  
  
or  
  
{}  
    'action': 'action_name',  
    'object': 'target_obj_name1,target_obj_name2',  
}
```

The action\_name must be one of the following:

```
LEFT_GRASP # the robot grasps the object with its left hand, to execute the action, the robot's left hand must be empty, e.g. {'action': 'LEFT_GRASP', 'object': 'apple_0'}.
RIGHT_GRASP # the robot grasps the object with its right hand, to execute the action, the robot's right hand must be empty, e.g. {'action': 'RIGHT_GRASP', 'object': 'apple_0'}.
LEFT_PLACE_ONTOP # the robot places the object in its left hand on top of the target object and release the object in its left hand, e.g. {'action': 'LEFT_PLACE_ONTOP', 'object': 'table_1'}.
RIGHT_PLACE_ONTOP # the robot places the object in its right hand on top of the target object and release the object in its left hand, e.g. {'action': 'RIGHT_PLACE_ONTOP', 'object': 'table_1'}.
LEFT_PLACE_INSIDE # the robot places the object in its left hand inside the target object and release the object in its left hand, to execute the action, the robot's left hand must hold an object, and the target object can't be closed e.g. {'action': 'LEFT_PLACE_INSIDE', 'object': 'fridge_1'}.
RIGHT_PLACE_INSIDE # the robot places the object in its right hand inside the target object and release the object in its left hand, to execute the action, the robot's right hand must hold an object, and the target object can't be closed, e.g. {'action': 'RIGHT_PLACE_INSIDE', 'object': 'fridge_1'}.
RIGHT_RELEASE # the robot directly releases the object in its right hand, to execute the action, the robot's left hand must hold an object, e.g. {'action': 'RIGHT_RELEASE', 'object': 'apple_0'}.
LEFT_RELEASE # the robot directly releases the object in its left hand, to execute the action, the robot's right hand must hold an object, e.g. {'action': 'LEFT_RELEASE', 'object': 'apple_0'}.
OPEN # the robot opens the target object, to execute the action, the target object should be openable and closed, also, toggle off the target object first if want to open it, e.g. {'action': 'OPEN', 'object': 'fridge_1'}.
CLOSE # the robot closes the target object, to execute the action, the target object should be openable and open, e.g. {'action': 'CLOSE', 'object': 'fridge_1'}.
COOK # the robot cooks the target object, to execute the action, the target object should be put in a pan, e.g. {'action': 'COOK', 'object': 'apple_0'}.
CLEAN # the robot cleans the target object, to execute the action, the robot should have a cleaning tool such as rag, the cleaning tool should be soaked if possible, or the target object should be put into a toggled on cleaner like a sink or a dishwasher, e.g. {'action': 'CLEAN', 'object': 'window_0'}.
FREEZE # the robot freezes the target object e.g. {'action': 'FREEZE', 'object': 'apple_0'}.
UNFREEZE # the robot unfreezes the target object, e.g. {'action': 'UNFREEZE', 'object': 'apple_0'}.
SLICE # the robot slices the target object, to execute the action, the robot should have a knife in hand, e.g. {'action': 'SLICE', 'object': 'apple_0'}.
SOAK # the robot soaks the target object, to execute the action, the target object must be put in a toggled on sink, e.g. {'action': 'SOAK', 'object': 'rag_0'}.
DRY # the robot dries the target object, e.g. {'action': 'DRY', 'object': 'rag_0'}.
```

```

TOGGLE_ON # the robot toggles on the target object, to execute the action, the target
          object must be closed if the target object is openable and open e.g. {{'action':
          ': 'TOGGLE_ON', 'object': 'light_0'}}.
TOGGLE_OFF # the robot toggles off the target object, e.g. {{'action': 'TOGGLE_OFF',
          'object': 'light_0'}}.
LEFT_PLACE_NEXTTO # the robot places the object in its left hand next to the target
                  object and release the object in its left hand, e.g. {{'action': '
                  LEFT_PLACE_NEXTTO', 'object': 'table_1'}}.
RIGHT_PLACE_NEXTTO # the robot places the object in its right hand next to the target
                  object and release the object in its right hand, e.g. {{'action': '
                  RIGHT_PLACE_NEXTTO', 'object': 'table_1'}}.
LEFT_TRANSFER_CONTENTS_INSIDE # the robot transfers the contents in the object in its
                  left hand inside the target object, e.g. {{'action': '
                  LEFT_TRANSFER_CONTENTS_INSIDE', 'object': 'bow_1'}}.
RIGHT_TRANSFER_CONTENTS_INSIDE # the robot transfers the contents in the object in its
                  right hand inside the target object, e.g. {{'action': '
                  RIGHT_TRANSFER_CONTENTS_INSIDE', 'object': 'bow_1'}}.
LEFT_TRANSFER_CONTENTS_ONTOP # the robot transfers the contents in the object in its
                  left hand on top of the target object, e.g. {{'action': '
                  LEFT_TRANSFER_CONTENTS_ONTOP', 'object': 'table_1'}}.
RIGHT_TRANSFER_CONTENTS_ONTOP # the robot transfers the contents in the object in its
                  right hand on top of the target object, e.g. {{'action': '
                  RIGHT_TRANSFER_CONTENTS_ONTOP', 'object': 'table_1'}}.
LEFT_PLACE_NEXTTO_ONTOP # the robot places the object in its left hand next to target
                  object 1 and on top of the target object 2 and release the object in its left
                  hand, e.g. {{'action': 'LEFT_PLACE_NEXTTO_ONTOP', 'object': 'window_0, table_1
                  '}}.
RIGHT_PLACE_NEXTTO_ONTOP # the robot places the object in its right hand next to
                  object 1 and on top of the target object 2 and release the object in its right
                  hand, e.g. {{'action': 'RIGHT_PLACE_NEXTTO_ONTOP', 'object': 'window_0, table_1
                  '}}.
LEFT_PLACE_UNDER # the robot places the object in its left hand under the target
                  object and release the object in its left hand, e.g. {{'action': '
                  LEFT_PLACE_UNDER', 'object': 'table_1'}}.
RIGHT_PLACE_UNDER # the robot places the object in its right hand under the target
                  object and release the object in its right hand, e.g. {{'action': '
                  RIGHT_PLACE_UNDER', 'object': 'table_1'}}.

Format of the interactable objects:
Interactable object will contain multiple lines, each line is a dictionary with the
following format:
{{{
  'name': 'object_name',
  'category': 'object_category'
}}}
object_name is the name of the object, which you must use in the action command,
object_category is the category of the object, which provides a hint for you in
interpreting initial and goal conditions.

Please pay special attention:
1. The robot can only hold one object in each hand.
2. Action name must be one of the above action names, and the object name must be one
   of the object names listed in the interactable objects.
3. All PLACE actions will release the object in the robot's hand, you don't need to
   explicitly RELEASE the object after the PLACE action.
4. For LEFT_PLACE_NEXTTO_ONTOP and RIGHT_PLACE_NEXTTO_ONTOP, the action command are
   in the format of {{'action': 'action_name', 'object': 'obj_name1, obj_name2'}}}
5. If you want to perform an action to an target object, you must make sure the
   target object is not inside a closed object.
6. For actions like OPEN, CLOSE, SLICE, COOK, CLEAN, SOAK, DRY, FREEZE, UNFREEZE,
   TOGGLE_ON, TOGGLE_OFF, at least one of the robot's hands must be empty, and the
   target object must have the corresponding property like they're openable,
   toggleable, etc.
7. For PLACE actions and RELEASE actions, the robot must hold an object in the
   corresponding hand.
8. Before slicing an object, the robot can only interact with the object (e.g.
   peach_0), after slicing the object, the robot can only interact with the sliced
   object (e.g. peach_0_part_0).

```

Examples: after# is the explanation.

```

Example 1:
Input:
initial environment state:
['stained', 'sink_7']
['stained', 'bathtub_4']
['not', 'soaked', 'rag_0']
['onfloor', 'rag_0', 'room_floor_bathroom_0']

```

```

['inside', 'rag_0', 'cabinet_1']
['not', 'open', 'cabinet_1']

target environment state:
['not', 'stained', 'bathtub_4']
['not', 'stained', 'sink_7']
['and', 'soaked', 'rag_0', 'inside', 'rag_0', 'bucket_0']

interactable objects:
[{'name': 'sink_7', 'category': 'sink.n.01'}
 {'name': 'bathtub_4', 'category': 'bathtub.n.01'}
 {'name': 'bucket_0', 'category': 'bucket.n.01'}
 {"name": "rag_0", "category": "rag.n.01"}
 {"name": "cabinet_1", "category": "cabinet.n.01"}]

Please output the list of action commands (in the given format) so that after the
robot executes the action commands sequentially, the current environment state
will change to target environment state. Usually, the robot needs to execute
multiple action commands consecutively to achieve final state. Please output
multiple action commands rather than just one. Only output the list of action
commands with nothing else.

Output:
[
    {
        "action": "OPEN",
        "object": "cabinet_1"
    }, # you want to get the rag_0 from cabinet_1, should open it first
    {
        "action": "RIGHT_GRASP",
        "object": "rag_0"
    }, # you want to clean the sink_7 and bathtub_4, you found them stained, so you
    need to soak the rag_0 first
    {
        "action": "RIGHT_PLACE_INSIDE",
        "object": "sink_7"
    }, # to soak the rag_0, you need to place it inside the sink_7
    {
        "action": "TOGGLE_ON",
        "object": "sink_7"
    }, # to soak the rag_0, you need to toggle on the sink_7
    {
        "action": "SOAK",
        "object": "rag_0"
    }, # now you can soak the rag_0
    {
        "action": "TOGGLE_OFF",
        "object": "sink_7"
    }, # after soaking the rag_0, you need to toggle off the sink_7
    {
        "action": "LEFT_GRASP",
        "object": "rag_0"
    }, # now you can grasp soaked rag_0 to clean stain
    {
        "action": "CLEAN",
        "object": "sink_7"
    }, # now you clean the sink_7
    {
        "action": "CLEAN",
        "object": "bathtub_4"
    }, # now you clean the bathtub_4
    {
        "action": "LEFT_PLACE_INSIDE",
        "object": "bucket_0"
    } # after cleaning the sink_7, you need to place the rag_0 inside the bucket_0
]

Your task:
Input:
initial environment state:
{init_state}

target environment state:
{target_state}

interactable objects:

```

2112

2113

2114

#### 2115 L.4 Prompt of Transition Modeling

2116 We design the prompt template of transition modeling as below:

**Prompt of Transition Modeling**

```
# Virtualhome Transition Modeling Prompt Template

SYSTEM_PROMPT = "You are a software engineer who will be writing action definitions
for a household robot in the PDDL planning language given the problem file and
predicates in domain file. For this task, please only output a parsable json
string inside brackets. Please start your answer with { and end your answer with
}. Don't include any notes or explanations with the output json string."

USER_PROMPT = "The following is predicates defined in this domain file. Pay attention
to the types for each predicate.
(define (domain virtualhome)
(:requirements :typing)
;; types in virtualhome domain
(:types
    object character ; Define 'object' and 'character' as types
)

;; Predicates defined on this domain. Note the types for each predicate.
(:predicates
    (closed ?obj - object) ; obj is closed
    (open ?obj - object) ; obj is open
    (on ?obj - object) ; obj is turned on, or it is activated
    (off ?obj - object) ; obj is turned off, or it is deactivated
    (plugged_in ?obj - object) ; obj is plugged in
    (plugged_out ?obj - object) ; obj is unplugged
    (sitting ?char - character) ; char is sitting, and this represents a state
of a character
    (lying ?char - character) ; char is lying
    (clean ?obj - object) ; obj is clean
    (dirty ?obj - object) ; obj is dirty
    (obj_ontop ?obj1 ?obj2 - object) ; obj1 is on top of obj2
    (ontop ?char - character ?obj - object) ; char is on obj
    (on_char ?obj - object ?char - character) ; obj is on char
    (inside_room ?obj ?room - object) ; obj is inside room
    (obj_inside ?obj1 ?obj2 - object) ; obj1 is inside obj2
    (inside ?char - character ?obj - object) ; char is inside obj
    (obj_next_to ?obj1 ?obj2 - object) ; obj1 is close to or next to obj2
    (next_to ?char - character ?obj - object) ; char is close to or next to obj
    (between ?obj1 ?obj2 ?obj3 - object) ; obj1 is between obj2 and obj3
    (facing ?char - character ?obj - object) ; char is facing obj
    (holds_rh ?char - character ?obj - object) ; char is holding obj with right
hand
    (holds_lh ?char - character ?obj - object) ; char is holding obj with left
hand
    (grabbable ?obj - object) ; obj can be grabbed
    (cuttable ?obj - object) ; obj can be cut
    (can_open ?obj - object) ; obj can be opened
    (readable ?obj - object) ; obj can be read
    (has_paper ?obj - object) ; obj has paper
    (movable ?obj - object) ; obj is movable
    (pourable ?obj - object) ; obj can be poured from
    (cream ?obj - object) ; obj is cream
    (has_switch ?obj - object) ; obj has a switch
```

2117

```

(lookable ?obj - object) ; obj can be looked at
(has_plug ?obj - object) ; obj has a plug
(drinkable ?obj - object) ; obj is drinkable
(body_part ?obj - object) ; obj is a body part
(recipient ?obj - object) ; obj is a recipient
(containers ?obj - object) ; obj is a container
(cover_object ?obj - object) ; obj is a cover object
(surfaces ?obj - object) ; obj has surfaces
(sittable ?obj - object) ; obj can be sat on
(lieable ?obj - object) ; obj can be lied on
(person ?obj - object) ; obj is a person
(hangable ?obj - object) ; obj can be hanged
(clothes ?obj - object) ; obj is clothes
(eatable ?obj - object) ; obj is eatable
)
;;
Actions to be predicted
)



Objective: Given the problem file of pddl, which defines objects in the task (:objects), initial conditions (:init) and goal conditions (:goal), write the body of PDDL actions (:precondition and :effect) given specific action names and parameters, so that after executing the actions in some order, the goal conditions can be reached from initial conditions.



Each PDDL action definition consists of four main components: action name, parameters , precondition, and effect. Here is the general format to follow:



```
(:action [action name]
  :parameters ([action parameters])
  :precondition ([action precondition])
  :effect ([action effect])
)
```



The :parameters is the list of variables on which the action operates. It lists variable names and variable types.



The :precondition is a first-order logic sentence specifying preconditions for an action. The precondition consists of predicates and 4 possible logical operators : or, and, not, exists!



1. The precondition should be structured in Disjunctive Normal Form (DNF), meaning an OR of ANDs.
2. The not operator should only be used within these conjunctions. For example, (or (and (predicate1 ?x) (predicate2 ?y)) (and (predicate3 ?x)))
3. Exists operator is followed by two parts, variable and body. It follows the format : exists (?x - variable type) (predicate1 ?x), which means there exists an object ?x of certain variable type, that predicate1 ?x satisfies.



The :effect lists the changes which the action imposes on the current state. The precondition consists of predicates and 6 possible logical operators: or, and, not, exists, when, forall.



1. The effects should generally be several effects connected by AND operators.
2. For each effect, if it is a conditional effect, use WHEN to check the conditions. The semantics of (when [condition] [effect]) are as follows: If [condition] is true before the action, then [effect] occurs afterwards.
3. If it is not a conditional effect, use predicates directly.
4. The NOT operator is used to negate a predicate, signifying that the condition will not hold after the action is executed.
5. Forall operator is followed by two parts, variable and body. It follows the format : forall (?x - variable type) (predicate1 ?x), which means there for all objects ?x of certain variable type, that predicate1 ?x satisfies.
6. An example of effect is (and (when (predicate1 ?x) (not (predicate2 ?y))) (predicate3 ?x))



Formally, the preconditions and effects are all clauses <Clause>.



```
<Clause> ::= (predicate ?x)
<Clause> ::= (and <Clause1> <Clause2> ...)
<Clause> ::= (or <Clause1> <Clause2> ...)
<Clause> ::= (not <Clause>)
<Clause> ::= (when <Clause1> <Clause2>)
<Clause> ::= (exists (?x - object type) <Clause>)
<Clause> ::= (forall (?x - object type) <Clause>)
```



In any case, the occurrence of a predicate should agree with its declaration in terms of number and types of arguments defined in DOMAIN FILE at the beginning.



Here is an example of the input problem file and unfinished action. Observe carefully how to think step by step to write the action body of hang_up_clothes:



Input:  
Problem file:


```

```

(define (problem hang-clothes-problem)
  (:domain household)
  (:objects
    character - character
    shirt - object
    hanger - object
  ) ; This section declares the instances needed for the problem: character is an
    instance of a character; shirt is an instance of an object classified as clothes
    ; hanger is an object that is suitable for hanging clothes.
  (:init
    (clothes shirt)
    (hangable hanger)
    (holds_rh alice shirt)
    (next_to alice hanger)
  ) ; This section declares the initial conditions. (clothes shirt) and (hangable
    hanger) tells the properties of objects; (holds_rh alice shirt) indicates that
    Alice is holding the shirt in her right hand; (next_to alice hanger) means Alice
    is next to the hanger, ready to hang the shirt.
  (:goal
    (and
      (ontop shirt hanger)
    )
  ) ; This section declares the goal. (ontop shirt hanger) is the goal, where the
    shirt should end up hanging on the hanger.
)
Action to be finished:
(:action hang_up_clothes
  :parameters (?char - character ?clothes - object ?hang_obj - object)
  :precondition ()
  :effect ())
)

Example output:
Given the objects in the problem file, and what typically needs to be true to perform
an action like hanging up clothes: 1. clothes must indeed be a type of clothing
. 2. hang_obj should be something on which clothes can be hung (hangable). 3.
char should be holding the clothes, either in the right or left hand. 4. char
needs to be next to the hanging object to hang the clothes. Besides, we need to
write preconditions in Disjunctive Normal Form.
These insights guide us to write:
:precondition (or
  (and
    (clothes ?clothes) ; the object must be a piece of clothing
    (hangable ?hang_obj) ; the target must be an object suitable for
    hanging clothes
    (holds_rh ?char ?clothes) ; character is holding clothes in the
    right hand
    (next_to ?char ?hang_obj) ; character is next to the hanging
    object
  )
  (and
    (clothes ?clothes) ; the object must be a piece of clothing
    (hangable ?hang_obj) ; the target must be an object suitable for
    hanging clothes
    (holds_lh ?char ?clothes) ; character is holding clothes in the
    left hand
    (next_to ?char ?hang_obj) ; character is next to the hanging
    object
  )
)
Effects describe how the world state changes due to the action. After hanging up
clothes, you'd expect: 1. char is no longer holding the clothes. 2. clothes is
now on the hang_obj.
These expectations convert into effects:
:effect (and
  (when (holds_rh ?char ?clothes) (not (holds_rh ?char ?clothes))) ; if
  clothes are held in the right hand, they are no longer held
  (when (holds_lh ?char ?clothes) (not (holds_lh ?char ?clothes))) ; if
  clothes are held in the left hand, they are no longer held
  (ontop ?clothes ?hang_obj) ; clothes are now hanging on the object
)

Combining these parts, the complete hang_up_clothes action becomes:
(:action hang_up_clothes
  :parameters (?char - character ?clothes - object ?hang_obj - object)
  :precondition (or
    (and
      (clothes ?clothes)
      (hangable ?hang_obj)
    )
  )
  :effect (and
    (when (holds_rh ?char ?clothes) (not (holds_rh ?char ?clothes))) ; if
    clothes are held in the right hand, they are no longer held
    (when (holds_lh ?char ?clothes) (not (holds_lh ?char ?clothes))) ; if
    clothes are held in the left hand, they are no longer held
    (ontop ?clothes ?hang_obj) ; clothes are now hanging on the object
  )
)
)

```

```

        (holds_rh ?char ?clothes)
        (next_to ?char ?hang_obj)
    )
    (and
        (clothes ?clothes)
        (hangable ?hang_obj)
        (holds_lh ?char ?clothes)
        (next_to ?char ?hang_obj)
    )
)
:effect (and
    (when (holds_rh ?char ?clothes) (not (holds_rh ?char ?clothes)))
    (when (holds_lh ?char ?clothes) (not (holds_lh ?char ?clothes)))
    (ontop ?clothes ?hang_obj)
)
)

```

Above is a good example of given predicates in domain file, problem file, action names and parameters, how to reason step by step and write the action body in PDDL. Pay attention to the usage of different connectives and their underlying logic.

Here are some other commonly used actions and their PDDL definition:

```

(:action put_to
:parameters (?char - character ?obj - object ?dest - object)
:precondition (or
    (and
        (hold_lh ?obj) ; The character should hold either with left hand or
        right hand
        (next_to ?char ?dest) ; The character should be close to destination
    )
    (and
        (hold_rh ?obj) ; The character should hold either with left hand or
        right hand
        (next_to ?char ?dest) ; The character should be close to destination
    )
)
:effect (obj_ontop ?obj ?dest) ; The object is now on the destination
)

```

This case illustrates the use of OR to include all possible preconditions of an action.

```

(:action pick_and_place
:parameters (?char - character ?obj - object ?dest - object)
:precondition (and
    (grabbable ?obj) ; The object must be grabbable
    (next_to ?char ?obj) ; The character must be next to the object
    (not (obj_ontop ?obj ?dest)) ; Ensure the object is not already on the
    destination
)
:effect (and
    (obj_ontop ?obj ?dest) ; The object is now on the destination
    (next_to ?char ?dest) ; The character is now next to the destination
)
)

```

This case illustrates a plain case with only AND operator.

```

(:action bow
:parameters (?char - character ?target - character)
:precondition (and
    (next_to ?char ?target) ; The character must be next to the target to
    perform the bow
)
:effect ()
)

```

This case illustrates the action can have no effect (or no precondition.)

hint:

1. Don't enforce the use of WHEN everywhere.
2. You MUST only use predicates and object types exactly as they appear in the domain file at the beginning.
3. Use and only use the arguments provided in :parameters for each action. Don't propose additional arguments, unless you are using exists or forall.
4. It is possible that action has no precondition or effect.

5. The KEY of the task is to ensure after executing your proposed actions in some order, the initial state (:init) in problem file can reach the goals (:goal)!!! Pay attention to the initial state and final goals in problem file.

6. Preconditions and effects are <Clause> defined above. When there is only one predicate, do not use logic connectives.

For actions to be finished, write their preconditions and effects, and return in standard PDDL format:

```
(:action [action name]
  :parameters ([action parameters])
  :precondition ([action precondition])
  :effect ([action effect])
)
```

Concatenate all actions PDDL string into a single string. Output in json format where key is 'output' and value is your output string: {'output': YOUR OUTPUT STRING}

Input:  
<problem\_file>  
<action\_handlers>

Output:"

2121

2122

2123

## 2124 L.5 Prompt of Environment Representation

2125 We input the object-centric representation as the abstraction of the environment as below:

### Prompt of Environment Abstraction

```
# Behavior Environment Prompt Template
AVAILABLE_STATES = '''The state is represented as a first-order predicate, which is a tuple of a predicate name and its arguments. Its formal definition looks like this "<PredicateName>(Params)", where <PredicateName> is the state name and each param should be ended with an id. An example is "inside(coin.1, jar.1)". Below is a list of available states and their descriptions.
| State Name | Arguments | Description |
| --- | --- | --- |
| inside | (obj1.id, obj2.id) | obj1 is inside obj2. If we have state inside(A, B), and you want to take A out of B while B is openable and stayed at "not open" state, please open B first. Also, inside(obj1, agent) is invalid.|
| ontop | (obj1.id, obj2.id) | obj1 is on top of obj2 |
| nextto | (obj1.id, obj2.id) | obj1 is next to obj2 |
| under | (obj1.id, obj2.id) | obj1 is under obj2 |
| onfloor | (obj1.id, floor2.id) | obj1 is on the floor2 |
| touching | (obj1.id, obj2.id) | obj1 is touching or next to obj2 |
| cooked | (obj1.id) | obj1 is cooked |
| burnt | (obj1.id) | obj1 is burnt |
| dusty | (obj1.id) | obj1 is dusty. If want to change dusty(obj1.id) to "not dusty(obj1.id)", there are two ways to do it, depending on task conditions. Here , all objects other than obj1 are types but not instances: 1. [inside(obj1.id, dishwasher or sink), toggledon(dishwasher or sink)] 2. holding other cleannning tool |
| frozen | (obj1.id) | obj1 is frozen |
| open | (obj1.id) | obj1 is open |
| sliced | (obj1.id) | obj1 is sliced. If want to change "not sliced(obj1.id)" to "sliced(obj1.id)", one must have a slicer. |
| soaked | (obj1.id) | obj1 is soaked |
| stained | (obj1.id) | obj1 is stained. If want to change stained(obj1.id) to "not stained(obj1.id)", there are three ways to do it, depending on task conditions. Here, all objects other than obj1 are types but not instances: 1. [inside(obj1.id, sink), toggledon(sink)] 2. [soaked(cleaner)] 3. holding detergent. |
| toggledon | (obj1.id) | obj1 is toggled on |
| holds_rh | (obj1.id) | obj1 is in the right hand of the robot |
| holds_lh | (obj1.id) | obj1 is in the left hand of the robot |'''
```

```
AVAILABLE_CONNECTIVES = '''The connectives are used to satisfy the complex conditions . They are used to combine the state predicates.
| Connective Name | Arguments | Description |
| --- | --- | --- |
| and | exp1 and exp2 | evaluates to true if both exp1 and exp2 are true |
```

2126

```

| or | exp1 or exp2 | evaluates to true if either exp1 or exp2 is true |
| not | not exp | evaluates to true if exp is false |
| forall | forall(x, exp) | evaluates to true if exp is true for all x |
| exists | exists(x, exp) | evaluates to true if exp is true for at least one x |
| forpairs | forpairs(x, y, exp) | evaluates to true if exp is true for all pairs
of x and y. For example, forpairs(watch, basket, inside(watch, basket)) means
that for each watch and basket, the watch is inside the basket. |
| forn | forn(n, x, exp) | evaluates to true if exp is true for exactly n times
for x. For example, forn(2, jar_n_01, (not open(jar_n_01))) means that there are
exactly two jars that are not open. |
| fornpairs | fornpairs(n, x, y, exp) | evaluates to true if exp is true for
exactly n times for pairs of x and y. For example, fornpairs(2, watch, basket,
inside(watch, basket)) means that there are exactly two watches inside the
basket. |''''

EXAMPLE_RELEVANT_OBJECTS = '''## Relevant objects in this scene
{'name': 'strawberry.0', 'category': 'strawberry_n_01'}
{'name': 'fridge.97', 'category': 'electric_refrigerator_n_01'}
{'name': 'peach.0', 'category': 'peach_n_03'}
{'name': 'countertop.84', 'category': 'countertop_n_01'}
{'name': 'jar.0', 'category': 'jar_n_01'}
{'name': 'jar.1', 'category': 'jar_n_01'}
{'name': 'carving_knife.0', 'category': 'carving_knife_n_01'}
{'name': 'bottom_cabinet_no_top.80', 'category': 'cabinet_n_01'}
{'name': 'room_floor_kitchen.0', 'category': 'floor_n_01'''}

INITIAL_STATES = '''inside(strawberry.0, fridge.97)
inside(peach.0, fridge.97)
not sliced(strawberry.0)
not sliced(peach.0)
ontop(jar.0, countertop.84)
ontop(jar.1, countertop.84)
ontop(carving_knife.0, countertop.84)
onfloor(agent_n_01.1, room_floor_kitchen.0)'''

GOAL_STATES = '''exists(jar_n_01, (inside(strawberry.0, jar_n_01) and (not inside(
peach.0, jar_n_01))))
exists(jar_n_01, (inside(peach.0, jar_n_01) and (not inside(strawberry.0,
jar_n_01))))
forall(jar_n_01, (not open(jar_n_01)))
sliced(strawberry.0)
sliced(peach.0)'''
```

2127

2128  
2129

Figure 49: Examples prompt of object-centric representation for the embodied environment.

## 2130 L.6 Prompt Analysis and Learned Lessons

- 2131 During our empirical experiments, we find many models struggle with outputting strictly formatted parsable  
2132 output that can be accepted by our embodied agent interface, let along the various intricate formats accepted  
2133 by different simulators. In addition, many models incline to output accompanying explanation along with their  
2134 answer despite explicit instructions forbidding such behavior.
- 2135 To address these problem, we design structured and easy-to-follow code output formats that consists of JSON  
2136 and Python list structures in place of long string output, leveraging popular LLMs' superior formal language  
2137 (code) generation ability over natural language.
- 2138 In order to standardize our evaluation, we use the same prompt for all LLMs evaluated, following the convention  
2139 established by [67].
- 2140 Below are some of the other LLM prompting strategies we used in designing our evaluation prompts. Based on  
2141 empirical qualitative experiments, prompts under these conditions generated the best results for our embodied  
2142 agent interface:

## Takeaways of Prompting Techniques

- **System Prompts** are preferred over user prompts for specifying output structure and overall output style as it is more strongly enforced.
- **Stop Sequences** are provided to prevent repetitive output and end-of-output explanations. This can improve format following and save token usage.
- **Shorter Strings** are encouraged to reduce grammar errors. For example, in the goal interpretation, we ask models to generate [“cond1”, “and”, “cond2”, “or”, “cond3”] instead of “cond1 and cond2 or cond3”.
- **Start / End Tokens** are specified to discourage models from generating additional explanation and self-repetition of answers / prompt instructions.

2143

## 2144 M Dataset Analysis

### 2145 M.1 Simulator Comparison and Selection

2146 The VirtualHome, BEHAVIOR, and AI2-THOR simulators each offer unique capabilities tailored to different  
2147 aspects of embodied agents. VirtualHome is ideal for long, complex household tasks, supporting a wide range of  
2148 scripted actions and detailed object states within home environments. BEHAVIOR provides a broad spectrum of  
2149 human activities with high goal complexity, realistic physics, and a large, dynamic state space, making it suitable  
2150 for testing generalization across diverse scenarios. AI2-THOR focuses on photorealistic indoor environments,  
2151 emphasizing detailed perception and manipulation tasks with rich real-time interactions and a substantial state  
2152 space.

2153 In detail, we compare their key differences in terms of task length, goal complexity, task scenarios, actions, and  
2154 state space:

#### 2155 Task Length and Goal Complexity

- 2156 • VirtualHome is designed to simulate longer sequences of everyday activities, involving complex  
2157 multi-step goals with many objects. It involves sequences of 10-30 high-level actions on average.
- 2158 • **BEHAVIOR** has the longest task length and most complex multi-step trajectories. Its range from 50  
2159 to over 1000 low-level actions in length.
- 2160 • **AI2-THOR** tasks are generally shorter, centered around navigation and basic object interactions with  
2161 around 10 steps.

#### 2162 Task Scenarios

- 2163 • **VirtualHome** specializes in simulating complete residential apartments/houses with furnished rooms  
2164 to enable realistic household activity scenarios such as like *making coffee* or *brush teeth*. It has has  
2165 around 20,000 unique household activity scenarios across 57 furnished residential environments.
- 2166 • **BEHAVIOR-100** provides 100 different scripted human behavior scenarios across various indoor  
2167 environments.
- 2168 • **AI2-THOR** focuses more on navigation tasks and low-level object interactions in diverse 3D indoor  
2169 scenes. It has has over 120 unique room configurations across 8 different room types like bathrooms,  
2170 living rooms, etc.

#### 2171 Action Space and State Space

- 2172 • **VirtualHome**: Represents activities as high-level action sequences like "<char0> [PutBack] <glass>  
2173 (1) <table>", enabling simulation of complex multi-step activities. It has around 300 unique high-level  
2174 action types.
- 2175 • **AI2-THOR**: Has a larger state space with many possible environment configurations, i.e., allowing  
2176 over 30 different low-level interactions like open, pick-up, toggle on/off across 200+ unique object  
2177 types. It has an estimated state space of over  $10^{25}$  possible environment configurations.
- 2178 • **BEHAVIOR**: Provides low-level scripted actions like object grasps and movements. It provides over  
2179 100 unique low-level action types like grasps, movements, etc.

#### 2180 Environments and Assets

- 2181 • **VirtualHome**: Has 57 fully furnished residential environments with over 3000 unique 3D object  
2182 assets.

- 2183 • **BEHAVIOR**: Provides 50 different indoor environments without physics simulation.  
 2184 • **AI2-THOR**: Contains over 120 unique room configurations with over 1000 3D object assets.

2185 Given our focus on embodied decision making, which emphasizes the ability to navigate complex goals and  
 2186 extended action sequences, we have selected VirtualHome and BEHAVIOR for our research. VirtualHome is  
 2187 particularly suited for simulating long, intricate household tasks with its extensive high-level action sequences  
 2188 and richly detailed environments. BEHAVIOR complements this by offering a diverse set of human activities  
 2189 with detailed, low-level scripted actions, enabling comprehensive testing of decision-making capabilities across  
 2190 various scenarios. Together, these platforms provide a robust framework for evaluating and enhancing the  
 2191 decision-making abilities of embodied agents in complex, goal-oriented tasks.

## 2192 M.2 Data Structure

2193 We define a data structure containing comprehensive annotation for the evaluation of four ability modules  
 2194 with different levels of planners. Each instance in the dataset represents a task goal, and each task contains the  
 2195 following data:

- 2196 1. Natural language task name  
 2197 2. Natural language task instruction  
 2198 3. Symbolic goal definition (including its LTL form)  
 2199 4. Symbolic action trajectory  
 2200 5. The transition models involved in the task

2201 As shown in Figure 6, we extend the RobotHow [2] data structure, where each task goal includes a natural  
 2202 language goal description, a VirtualHome action script, and the initial and final states of the scene. In addition,  
 2203 we offer precise symbolic goals for each task, which enhances the accuracy of evaluation by mitigating the impact  
 2204 of noisy final states. For instance, it is more precise to denote the final goal of the task “turn on light” as *on(light)*  
 2205 rather than a noisier version like *facing(agent, light) and on(light)*. Moreover, we provide transition model  
 2206 annotations, which require accurate annotation of all the logical constraints of preconditions and post-effects.

2207 For tasks in the BEHAVIOR environment, as shown in Figure 7, the dataset annotation focuses on the action  
 2208 sequences and transition models. The dataset also includes accompanying demo videos that showcase the  
 2209 execution of the ground truth action trajectories. For instance, the task “bottling fruit” includes a demo video  
 2210 showing the agent picking up peaches and placing them inside a jar, along with the corresponding action  
 2211 sequence and transition model annotations.

2212 The designed data structure is general and can support a systematic evaluation and usage of the data. It enables  
 2213 the exploration of different integration methods for various modules and can flexibly support downstream  
 2214 applications. By providing a comprehensive set of annotations, including natural language descriptions, symbolic  
 2215 goals, action trajectories, and transition models, researchers can investigate the interplay between different  
 2216 components of embodied AI systems and develop novel approaches for integrating them effectively. Also, it is  
 2217 not limited to specific simulators and can be expanded to other environments. It has the potential to become  
 2218 a standardized data format for embodied AI tasks, facilitating the comparison and benchmarking of different  
 2219 methods across various domains. The consistent representation of goals, actions, and transitions allows for a  
 2220 unified evaluation framework and promotes the development of generalizable and transferable approaches.

2221 The availability of natural language descriptions alongside symbolic representations enables the exploration of  
 2222 language grounding and the development of more intuitive human-robot interaction methods. For example, an  
 2223 agent can be trained to understand and execute natural language commands by grounding the language to the  
 2224 corresponding symbolic goals and actions.

## 2225 M.3 Goal Complexity Analysis

2226 Goals in BEHAVIOR may contain quantifiers, such as `forpairs( (?jar.n.01 - jar.n.01)`  
 2227 `(?apple.n.01 - apple.n.01)) (inside ?apple.n.01 ?jar.n.01)`, which are referred to  
 2228 as BDDL goals. These BDDL goals can be translated into multiple grounded goals. For exam-  
 2229 ple, the BDDL goal `forpairs( (?jar.n.01 - jar.n.01) (?apple.n.01 - apple.n.01))`  
 2230 `(inside ?apple.n.01 ?jar.n.01)` could be translated into the grounded goals `(and (inside`  
 2231 `apple.n.01_1 jar.n.01_1) (inside apple.n.01_2 jar.n.01_2))`. Different combina-  
 2232 tions of grounded goals can satisfy the same BDDL goal, which we refer to as goal options. For instance, `(and`  
 2233 `(inside apple.n.01_2 jar.n.01_1) (inside apple.n.01_1 jar.n.01_2))` also satis-  
 2234 fies the same BDDL goal. Generally, one BDDL goal may have multiple goal options, each consisting of  
 2235 several grounded goals.

2236 Figure 50 illustrates the distribution of the number of grounded goals and goal options in BEHAVIOR. The  
 2237 number of grounded goals is determined by sampling a goal option from the BDDL goal and counting the expres-  
 2238 sions within the brackets of the outermost connective *and*. For example, `(and (inside apple.n.01_2`

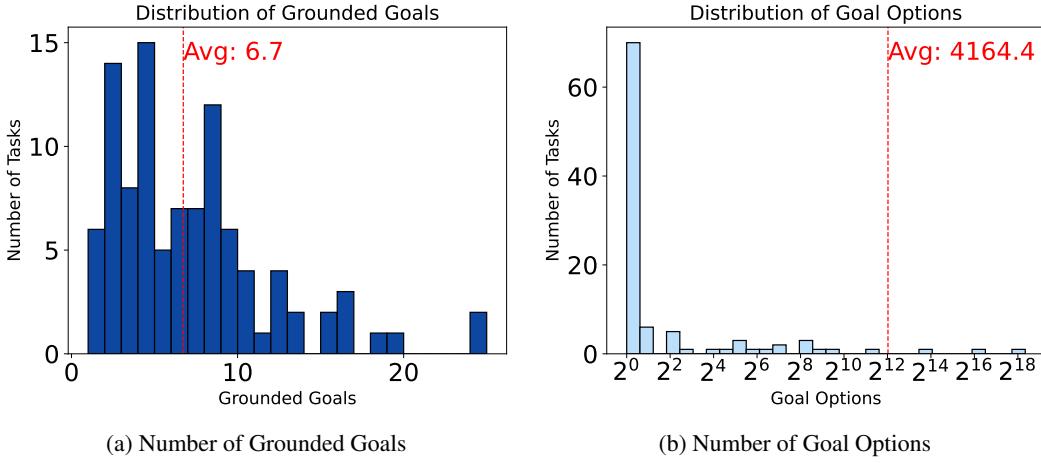


Figure 50: BEHAVIOR Goal Complexity

2239 jar.n.01\_1) (inside apple.n.01\_1 jar.n.01\_2)) contains two grounded goals. The number  
 2240 of goal options is calculated by the different possible combinations of grounded goals that satisfy the BDDL  
 2241 goal. On average, each task in BEHAVIOR has 6.7 grounded goals and 4164.4 goal options.

#### 2242 M.4 Data Statistics and Distribution

2243 To understand the datasets used for evaluating long-horizon decision making capabilities for complex goals, we  
 2244 provide detailed statistics and distribution information in the following sections.

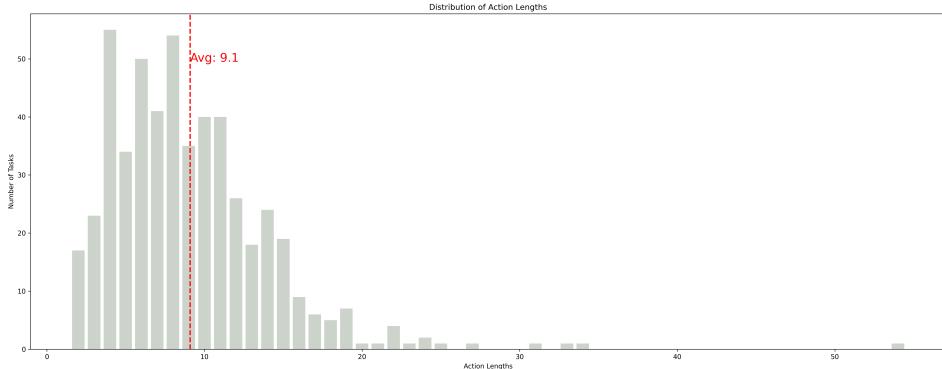


Figure 51: Action sequence length distribution in VirtualHome

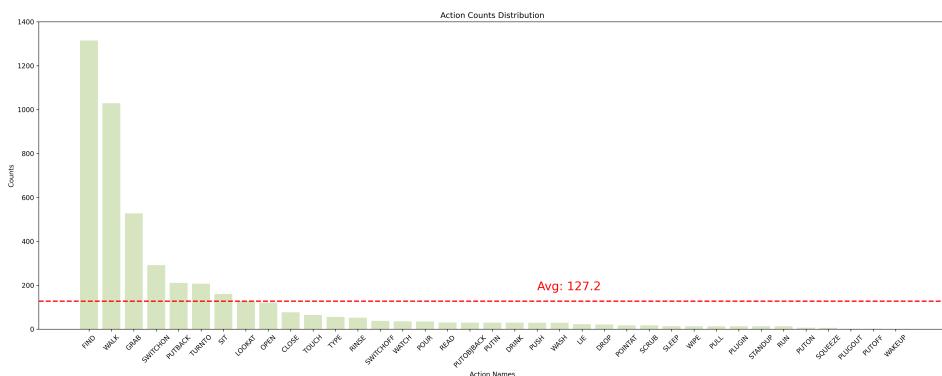


Figure 52: Action counts distribution in VirtualHome

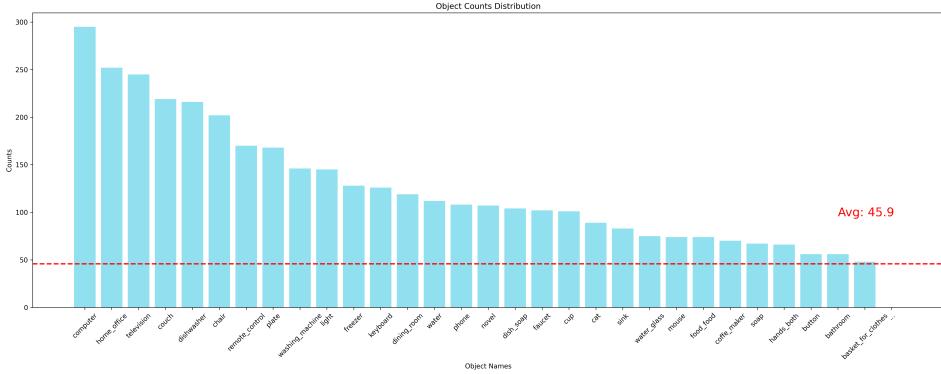


Figure 53: Object counts distribution in VirtualHome

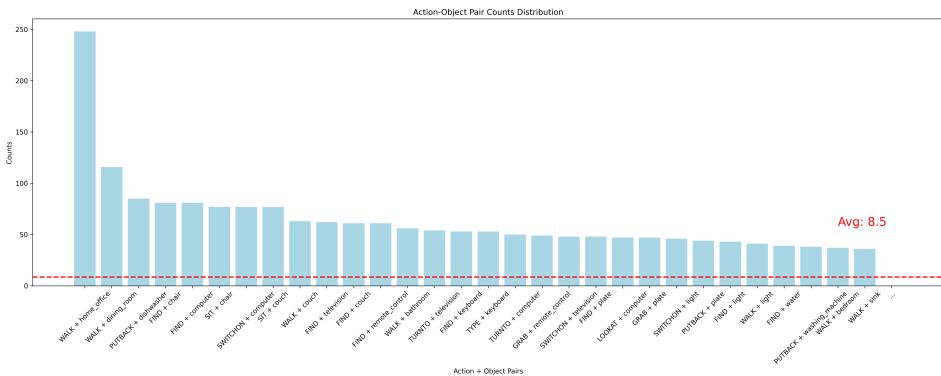


Figure 54: Action object pair counts distribution in VirtualHome

2245 **VirtualHome** Table 1 in the main paper presents the statistics of our annotated RobotHow dataset. We  
 2246 identified a total of 801 goals, comprising 340 state goals, 299 relation goals, and 162 action goals. For  
 2247 the provided task instructions (referred to as ‘trajectories’ in the table), the average length of action steps is  
 2248 approximately 8.76, ranging from 2 to 54 steps. Notably, 21 out of 26 task categories include instructions with  
 2249 action steps exceeding 10 in length, and one-third of the instructions have step lengths of more than 10. This  
 2250 indicates the complexity of our annotated data.

2251 **BEHAVIOR** As shown in Figure 55, the BEHAVIOR dataset has an average of 14.6 actions and 3.7  
 2252 goals per task. The action length is determined by the number of actions in the ground truth action se-  
 2253 quence. The number of goals is calculated by counting the expressions within the brackets of the outermost  
 2254 connective *and* in the BEHAVIOR goal, which is written in the Behavior Definition Description Language  
 2255 (BDDL) format. For example, for a goal expression such as `(and (contains ?hot_tub.n.02_1`  
 2256 `?chlorine.n.01_1) (filled ?hot_tub.n.02_1 ?water.n.06_1))`, we would consider it as  
 2257 having two goals: the first one is `(contains ?hot_tub.n.02_1 ?chlorine.n.01_1)`, and the  
 2258 second one is `(filled ?hot_tub.n.02_1 ?water.n.06_1)`.

## 2259 M.5 Task List

2260 Our dataset’s task selection criteria prioritize long-horizon tasks with complex goals, ensuring a challenging  
 2261 and comprehensive evaluation of embodied agents. For BEHAVIOR, we choose all the tasks with human  
 2262 demonstrations to ensure high-quality results. For VirtualHome, we select a typical scene and include all  
 2263 the tasks that can be executed within this scene. The combination of simulated tasks from VirtualHome and  
 2264 human-demonstrated tasks from BEHAVIOR allows to evaluate the generalization capabilities of LLMs across  
 2265 different domains and levels of complexity. The task lists in Table 20 and Table 21 offer a clear overview of the  
 2266 selected tasks to understand the scope and variety of the evaluation benchmark.

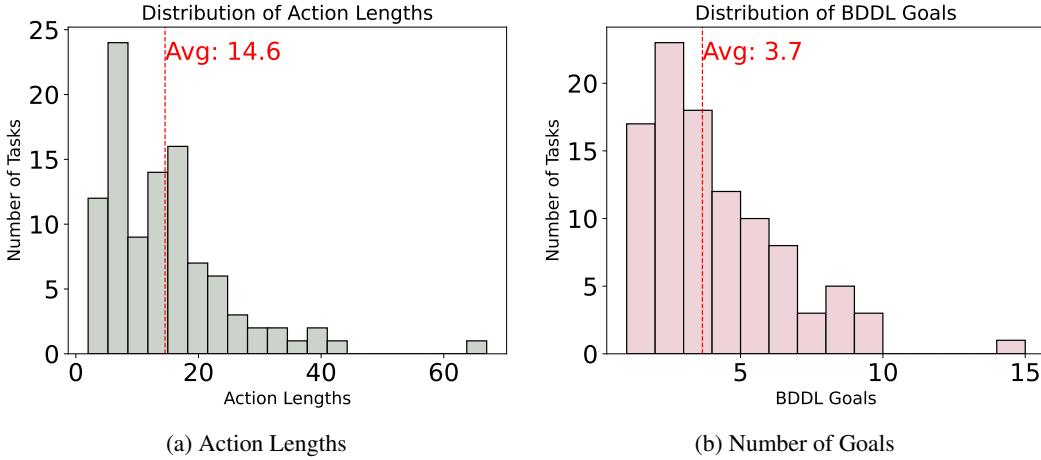


Figure 55: BEHAVIOR Task Complexity

Table 20: Task List on VirtualHome

Task Name	Task Number	Task Name	Task Number
Wash hands	22	Pet cat	27
Drink	28	Work	35
Turn on light	34	Pick up phone	28
Go to toilet	11	Get some water	1
Wash teeth	11	Read book	28
Make coffee	10	Listen to music	22
Wash dishes by hand	11	Relax on sofa	35
Wash clothes	21	Wash dishes with dishwasher	27
Browse internet	31	Put groceries in Fridge	23
Set up table	8	Write an email	21
Take shower	2	Watch TV	32
Cook some food	7	Change TV channel	29
Go to sleep	9		
Brush teeth	5		

2267 The inclusion of a typical scene in VirtualHome enables a focused evaluation of an agent’s performance within  
 2268 a specific environment. By providing all the executable tasks within this scene, researchers can thoroughly  
 2269 investigate the agent’s ability to navigate, interact with objects, and complete goals in a constrained setting.  
 2270 This setup facilitates the analysis of the agent’s behavior and helps identify strengths and weaknesses in its  
 2271 decision-making process. The selected tasks in VirtualHome cover a wide range of household activities, such as  
 2272 cooking, cleaning, and object manipulation. For instance, the task “cook some food” involves multiple steps,  
 2273 including gathering ingredients, using kitchen appliances, and setting the table. Another example is the task  
 2274 “wash dishes by hand”, which requires the agent to perform actions like picking up objects, wiping surfaces,  
 2275 and putting things back in their designated places. These tasks showcase the complexity and diversity of the  
 2276 scenarios.

Table 21: Task List on BEHAVIOR . Length is the action trajectory length.

Task Name	Length	Task Name	Length
assembling_gift_baskets	32	brushing_lint_off_clothing	13
boxing_books_up_for_storage	16	collecting_aluminum_cans	12
mopping_floors	15	preserving_food	20
re-shelving_library_books	8	polishing_silver	6
packing_boxes_for_household_move_or_trip	20	cleaning_freezer	15
installing_a_modem	3	cleaning_up_after_a_meal	31
putting_away_Christmas_decorations	17	setting_up_candles	12
cleaning_shoes	10	cleaning_sneakers	30
locking_every_window	4	washing_cars_or_other_vehicles	8
washing_dishes	11	putting_away_Halloween_decorations	15

Task Name	Frequency	Task Name	Frequency
organizing_boxes_in_garage	12	putting_up_Christmas_decorations_inside	18
cleaning_bathtub	8	clearing_the_table_after_dinner	14
loading_the_dishwasher	13	filling_an_Easter_basket	26
cleaning_the_pool	8	opening_packages	6
packing_bags_or_suitcase	15	cleaning_up_refrigerator	23
making_tea	13	picking_up_trash	10
polishing_furniture	4	organizing_school_stuff	15
cleaning_cupboards	24	installing_a_printer	3
packing_food_for_work	12	storing_the_groceries	23

2277 In the BEHAVIOR dataset, the chosen tasks are based on real-world human demonstrations, ensuring the  
 2278 authenticity and practicality of the tasks. These tasks encompass various domains, such as object manipulation,  
 2279 tool use, and goal-directed behaviors. For example, the task “bottling fruit” requires the agent to grasp peaches,  
 2280 cut them, and place them inside a jar, demonstrating object manipulation skills. The inclusion of these human-  
 2281 demonstrated tasks in BEHAVIOR provides valuable insights into real-world challenges and helps to evaluate  
 2282 the performance of embodied agents in more realistic settings.

## 2283 M.6 Task Categorization

2284 To have a better insight on LLM’s abilities on predicting transition modeling from different perspectives, we  
 2285 group all the predicates in PDDL into five categories: object affordance, object states, object orientation, spatial  
 2286 relations, and non-spatial relations.

2287 In our approach, we categorize programs based on the most dominant predicate category, employing Inverse  
 2288 Document Frequency (IDF) to measure the significance of predicates within each program. Given a set of  
 2289 programs  $P = \{p_1, p_2, \dots, p_n\}$  and a corresponding set of predicates  $T = \{t_1, t_2, \dots, t_m\}$ , each predicate  
 2290  $t_i$  is associated with a specific category  $C_j$  within a predefined set of categories  $C = \{c_1, c_2, \dots, c_k\}$ . Each  
 2291 action  $a_i$  corresponds to a set of predicates  $T'_{a_i} = \bigcup_{t_i \in \text{PDDL}(a_i)} t_i$  according to its PDDL definition  $\text{PDDL}(a_i)$ .  
 2292 Each program  $p_i$  can be solved by PDDL planner and has a corresponding action sequence  $A_i$ . The predicate  
 2293 set  $T_i$  is defined as the unification of the predicate sets of all actions involved in action sequence  $A_i$ . Formally,  
 2294  $T_i = \bigcup_{a_j \in A_i} T'_{a_j} = \{t_{i1}, t_{i2}, \dots, t_{il}\}$ , where  $l$  can vary between programs.

2295 The task categorization process begins with the calculation of the IDF for each predicate across the corpus of  
 2296 programs. The IDF for a predicate  $t$  is computed as follows:

$$\text{IDF}(t) = \log \left( \frac{|P|}{\text{df}(t)} \right)$$

2297 where  $|P|$  denotes the total number of programs and  $\text{df}(t)$  is the number of programs that contain the predicate  
 2298  $t$ . This metric quantifies the importance of a predicate relative to its frequency across all programs, thereby  
 2299 enhancing the significance of rarer predicates.

2300 Following the computation of the IDF values, each program  $p_i$  is scored for each category  $c_j$  by summing the  
 2301 IDF values of the predicates belonging to  $c_j$ :

$$S(c_j, p_i) = \sum_{t \in p_i, t \in c_j} \text{IDF}(t)$$

2302 Each program  $p_i$  is assigned to its top  $K$  scoring categories based on the accumulated IDF scores. This is  
 2303 formalized as follows:

$$\text{Categories}(p_i) = \text{top}_K \{(c_j, S(c_j, p_i)) \mid c_j \in C\}$$

2304 where  $\text{top}_K$  selects the set of  $K$  categories with the highest scores for each program.

2305 This scoring and categorization method emphasizes the content-specific importance of predicates within pro-  
 2306 grams, assigning categories based on the most informative predicate category present in each program. This  
 2307 approach is particularly effective in distinguishing the dominant themes of programs when predicates are  
 2308 unevenly distributed across categories, thus providing a robust basis for categorization in large-scale and diverse  
 2309 datasets.

2310 We categorize tasks with  $K = 2$ . Table 22 illustrates the number of tasks under each task category in VirtualHome  
 2311 and BEHAVIOR.

Table 22: Results on Task Categorization, showing task number in each task category.

	VirtualHome	BEHAVIOR
<b>Object States</b>	178	45
<b>Object Orientation</b>	167	-
<b>Object Affordance</b>	28	-
<b>Spatial Relations</b>	145	94
<b>Non-spatial Relations</b>	74	61

## 2312 N LTL Implementation Details

### 2313 N.1 Why LTL

2314 Our EMBODIED AGENT INTERFACE is built on top of the linear temporal logic (LTL) language. This is  
 2315 motivated by two critical desiderata of the interface. First, we need an expressive and compact language  
 2316 to describe task specifications. Classical choices such as first-order logic formulas on goal states or reward  
 2317 functions both have their limitations: goal state formulas only describe the requirements over the goal state  
 2318 but not any temporal ordering of how subgoals should be achieved. On the other hand, reward functions are  
 2319 general in specifying preferences over trajectories but they usually can not be represented in a compact way  
 2320 due to their numeric nature. Second, we need a unified interface between different modules of an embodi-  
 2321 ed agent system, such as the inputs and outputs for goal interpreters, subgoal generators, etc. For example,  
 2322 BEHAVIOR [1] uses BEHAVIOR Domain Definition Language (BDDL) to represent goals as a target state  
 2323 with logic constraints, such as “not(stained(fridge\_97)), forall tray.n.01-tray.n.01  
 2324 inside(tray.n.01,fridge\_97) not(stained(tray.n.01)), ...”. In contrast, Virtual-  
 2325 Home [2] describes task goals in natural language with a different focus, such as “*take everything out of*  
 2326 *the fridge, throw anything outdated...*”. Furthermore, different agents may follow different trajectories and  
 2327 have different criteria for achieving the same goal. For instance, BEHAVIOR focuses on state transitions to  
 2328 match final states with goals (“not stained(fridge)”), whereas VirtualHome evaluates execution success without  
 2329 verifying whether the goal states are satisfied. BEHAVIOR focuses on state transitions to match final states  
 2330 with goals (“not stained(fridge)”), while VirtualHome only considers execution success rates without checking  
 2331 whether goal states are satisfied. This leads to significant difference in goal interpretation and object state  
 2332 representation across environments.

2333 LTL provides an expressive and compact description language solution to these issues. At a high level, an LTL  
 2334 formula can describes basic state constraints (e.g., a subgoal should be achieved), action constraints (e.g., a  
 2335 particular action should be executed), and possible temporal orders and dependencies among them (e.g., all  
 2336 dishes should be cleaned before we start to cook). By combining temporal connectives such as “next” and  
 2337 propositional logic connectives, we can also flexibly describe alternative goals or subgoal sequences. Therefore,  
 2338 state goals, action sequences, pre- and post-conditions of actions, subgoal sequences, or even sets of candidate  
 2339 subgoal sequences, can all be expressed in LTL in a compact way. As a byproduct, using a single description  
 2340 language for all inputs and outputs enables us to design a unified evaluation metric to measure the prediction  
 2341 accuracy, by measuring the similarity between two LTL formulas, which is detailed in later sections.

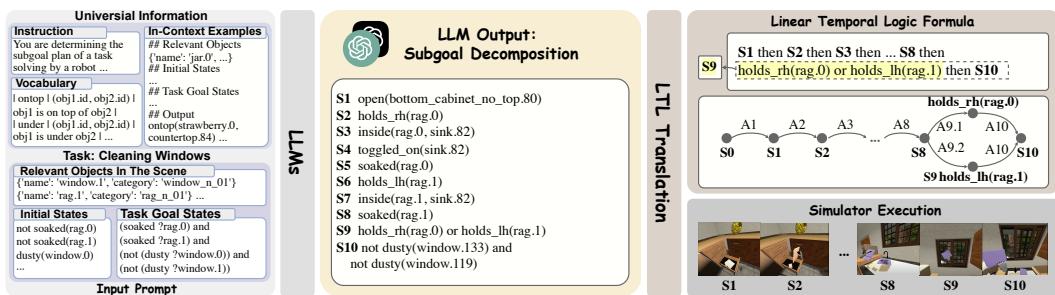


Figure 56: Pipeline of subgoal decomposition based on LTL in EMBODIED AGENT INTERFACE

2342 Figure 56 illustrates the complete process of subgoal decomposition using LTL in our EMBODIED AGENT  
 2343 INTERFACE. The process begins with describing the environment using LTL-like grammar. Once the prompt is  
 2344 crafted, a language model generates the corresponding output, which is then translated into a plain LTL formula.  
 2345 This formula is subsequently parsed into an LTL expression tree. Finally, a concrete subgoal path is sampled  
 2346 and converted into an action sequence. This action sequence is then executed in simulators to ensure two main  
 2347 criteria: (1) the subgoals are well-defined and executable, and (2) if executable, whether they meet the final state.

2348 **N.2 Syntax and Semantics of LTL Formulas**

2349 In EMBODIED AGENT INTERFACE, a state is represented as a tuple  $s = \langle U, F \rangle$ , where  $U$  is the universe of  
 2350 entities, assumed to be a fixed finite set.  $F$  is a set of relational Boolean features. Each feature  $f \in F$  can  
 2351 be viewed as a table where each entry is associated with a tuple of entities  $(o_1, \dots, o_k)$ . Each entry has the  
 2352 value of the feature in the state, and  $k$  is the arity of the feature. For example, the feature  $on(x, y)$  is a binary  
 2353 predicate. Actions can be viewed as primitive functions that take entities as inputs. For a physical robot,  
 2354 this corresponds to the available low-level controllers that our algorithm can interface with, such as moving and  
 2355 grasping.

2356 **LTL syntax.** Our EMBODIED AGENT INTERFACE uses a fragment of the full linear temporal logic (LTL)  
 2357 formalism on finite trajectories. In particular, we consider the following two types of atomic propositions. The  
 2358 arguments to these propositions can be either objects in the state (e.g., book1, cat1) or quantified variables (e.g.,  
 2359  $x$ ).

- 2360 (1) State propositions: Predicates that describe properties of object states and relations. For example,  
 2361  $ontop(book1, chair1)$ .  
 2362 (2) Action propositions: Predicates that denote actions. For example,  $touch(cat)$ .

2363 An LTL formula  $\phi$  is defined recursively as follows:

$$\phi ::= p \mid \neg\phi \mid \phi_1 \wedge \phi_2 \mid \phi_1 \vee \phi_2 \mid \phi_1 \Rightarrow \phi_2 \mid \forall x \phi(x) \mid \exists x \phi(x) \mid \exists^{=n} x \phi(x) \mid (\phi) \mid \phi_1 \text{ then } \phi_2$$

2364 where  $\phi_1$  and  $\phi_2$  are LTL formulas,  $p$  is an atomic proposition.  $\neg$  (negation),  $\wedge$  (and),  $\vee$  (or),  $\Rightarrow$  (implies)  
 2365 are logical connectives.  $\forall$ ,  $\exists$  and  $\exists^{=n}$  are quantifiers. Note that,  $\exists x$  means that there is at least one  $x$  such  
 2366 that  $\phi(x)$  is satisfied, whereas  $\exists^{=n} x$  means that there are exactly  $n$   $x$ 's such that  $\phi(x)$  is satisfied. **then** is a  
 2367 temporal connective, where  $\phi_1$  **then**  $\phi_2$  intuitively means  $\phi_1$  should happen before  $\phi_2$ <sup>8</sup>. Note that the operator  
 2368 **then** is a combination of the “next” and the “eventually” operator in standard LTL formalism, and we do not  
 2369 include “globally” and “until,” since the “then” operator is sufficient for describing all the task and input-output  
 2370 specifications in our system, although we can naturally extend our implementation to include them.

### LTL Grammar Definition

```
?start: stmt
primitive: VARNAME "(" [args] ")" # primitive format looks like varname(param)
object_name: VARNAME # object_name can be an object name (eg. pants)
           | VARNAMEWITHID # or object name with ID (eg. pants.1000)
args: object_name ("," object_name)* # definition of arguments

?stmt: then_stmt | primitive_stmt # a formula is a Boolean stmt
then_stmt: or_stmt ("then" or_stmt)* # connective priority order:
or_stmt: and_stmt ("or" and_stmt)* # then < or < and < not < forall = forn = exists
and_stmt: primitive_stmt ("and" primitive_stmt)*
primitive_stmt: "not" primitive_stmt -> not_stmt
              | primitive
              | "(" stmt ")"
              | "forall" VARNAME "." "(" stmt ")" -> forall_stmt
              | "forn" VARNAME "." "(" stmt ")" -> forn_stmt
              | "exists" VARNAME "." "(" stmt ")" -> exists_stmt

%import common.WS
%ignore WS
VARNAME: /[a-zA-Z_]\w*/
VARNAMEWITHID: /[a-zA-Z_]\w*\.[0-9]+/
```

2371

2372 **LTL semantics.** An LTL formula can be viewed as a classifier over trajectories semantically: we can evaluate  
 2373 an LTL formula  $\phi$  based on a state-action sequence. If the evaluation returns true, we say the state-action  
 2374 sequence satisfies  $\phi$ . This can be directly used to evaluate whether a generated action sequence satisfies the  
 2375 task specification. The task of a planner would be to take an LTL formula as its specification and generate a  
 2376 state-action sequence that satisfies the formula.

2377 Let a state-action trajectory  $T$  be  $[s_0, a_1, s_1, \dots, a_n, s_n]$ ,  $T_i = (s_i, a_i)$ , and  $U$  be the universe of entities  
 2378 in  $T$ . For a state-action pair, we can define the semantics of atomic propositions, logic connectives, and  
 2379 quantifiers. In particular, for atomic propositions  $p$ ,  $eval(p, (s_i, a_i))$  is true if  $p$  is satisfied in  $s_i$  (if  $p$  is  
 2380 a state predicate) or  $a_i = p$  (if  $p$  is an action predicate). All logic connectives ( $\neg$ ,  $\wedge$ ,  $\vee$ , and  $\Rightarrow$ ) and  
 2381 quantifiers ( $\forall$  and  $\exists$ ) follows their semantics in first-order logic. The for-n counting quantifier  $\exists^{=n}$  has the  
 2382 semantics that:  $eval(\exists^{=n} x. \phi(x), T_i) = \mathbb{1}[\sum_x eval(\phi(x), T_i) = n]$ , where  $\mathbb{1}[\cdot]$  is the indicator function. For  
 2383 compactness, if we apply a state-action formula  $\phi$  on a trajectory  $T$  instead of a concrete state-action pair  $T_i$ :  
 2384  $eval(\phi, T) = \exists k.eval(\phi, T_k)$ . That is,  $\phi$  is satisfied in at least one of the states in  $T$ .

<sup>8</sup>The priority of LTL operators from highest to lowest is () >  $\forall$  =  $\exists$  =  $\exists^{=n}$  >  $\neg$  >  $\wedge$  >  $\vee$  >  $\text{then}$ .

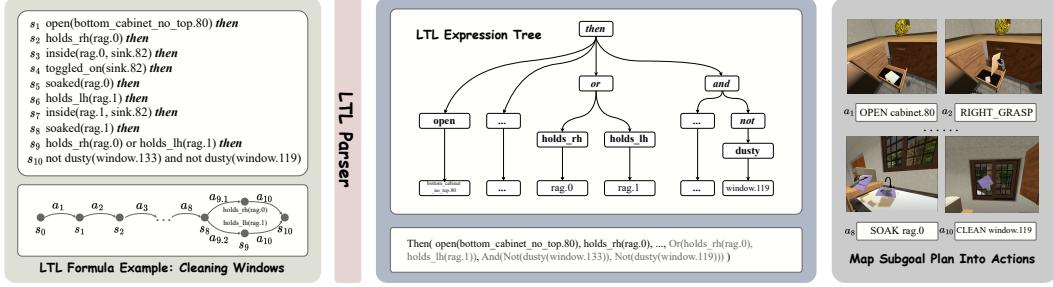


Figure 57: An example of LTL representation.

2385 The semantics of the operator **then** is defined as the following:

$$\text{eval}(\phi_1 \text{ then } \phi_2, T) = \exists k. \phi_1(T_{\leq k}) \wedge \phi_2(T_{>k}),$$

2386 where  $T_{\leq k}$  is the first  $k$  state-action pairs in  $T$  and  $T_{>k}$  is the suffix sequence after  $k$  steps. Intuitively, it means,  
2387 there exists a segmentation of the trajectory  $T$  such that  $\phi_1$  is satisfied in the first half while  $\phi_2$  is satisfied in the  
2388 second half.

2389 LTL formula will be parsed into an LTL expression tree before evaluation process, as demonstrated in Figure 57.  
2390 In order to evaluate the function  $\text{eval}(\phi, T)$  given the LTL formula and a state-action sequence, one needs to  
2391 recursively evaluate components in  $\phi$  based on their semantics. This is typically implemented with a dynamic  
2392 programming algorithm over LTL formulas and subsequences of  $T$ .

## 2393 O Extensive Related Work

### 2394 O.1 LLMs for Embodied Planning

2395 Existing works in embodied task and motion planning (TAMP) have used LLMs to perform varying tasks,  
2396 serving different ability modules defined within our embodied agent interface. Due to the page limit, we only  
2397 list a subset of such works and their categorization in the main paper. Here in Table 23 we provide a extended  
2398 list of such works with detailed categorization for your reference.

Table 23: Categorization of Existing Embodied Agent Planning Works’ Usage of Large Language Models: Each “LLMs” refers to the usage of LLMs to perform an ability module. For example, Ada [6] uses LLMs for Action Sequencing and Transition Modeling, while LLM+P [5] uses LLMs for Goal Interpretation.

Existing Work	Ref.	Goal Interpretation	Action Sequencing	Subgoal Decomposition	Transition Modeling
SayCan	[4]	LLMs	LLMs		
Ada	[6]	LLMs			LLMs
LLP+P	[5]	LLMs			
AutoTAMP	[15]		LLMs		LLMs
Code as Policies	[3]	LLMs	LLMs	LLMs	
Voyager	[31]	LLMs	LLMs		
Demo2Code	[10]	LLMs		LLMs	LLMs
LM as ZeroShot Planner	[62]		LLMs	LLMs	
SayPlan	[12]	LLMs	LLMs		LLMs
Text2Motion	[44]		LLMs		
LLMGROP	[42]	LLMs	LLMs		
REFLECT	[61]	LLMs	LLMs		
Generating Consistent PDDL Domains with LLMs	[28]	LLMs			LLMs
PlanSeqLearn	[59]		LLMs		
COWP	[66]	LLMs	LLMs		LLMs
HumanAssisted Robot Learning	Continual [60]		LLMs		

Continued on next page

Existing Works	Ref.	Goal Interpretation	Action Sequencing	Subgoal Decomposition	Transition Modeling
DECKARD	[30]	LLMs			LLMs
MOSAIC	[58]		LLMs		
Interactive Task Planning with Language Models	[11]		LLMs	LLMs	
RoCo	[57]		LLMs		
Cook2LTL	[38]	LLMs			
InnerMonologue	[18]		LLMs		
MLDT	[17]		LLMs		
Learning to Reason over Scene Graphs	[56]	LLMs	LLMs		LLMs
GRID	[16]	LLMs	LLMs		
LLMplanner	[29]	LLMs	LLMs		
DELTA	[14]		LLMs		
Look Before You Leap	[55]	LLMs	LLMs		
CAPE	[32]	LLMs	LLMs		
HERACLES	[39]		LLMs		
RoboTool	[13]		LLMs		LLMs
PROMST	[54]		LLMs		
LLM3	[27]	LLMs	LLMs		
Ghost in the Minecraft	[23]		LLMs		
PlanBench	[74]	LLMs	LLMs		
TaPA	[26]	LLMs	LLMs	LLMs	
ChatGPT Robot Control	[25]		LLMs		
LLM World Models for Planning	[24]	LLMs	LLMs		
DEPS	[72]	LLMs	LLMs		
Grounded Decoding	[22]		LLMs		
ProgPrompt	[33]	LLMs	LLMs		
DROC	[21]		LLMs		LLMs
LMPC	[53]	LLMs	LLMs		
GPTPDDL	[43]		LLMs		
RAP	[52]		LLMs		
LEAGUE++	[34]		LLMs		LLMs
CoPAL	[20]	LLMs	LLMs		
SayCanPay	[19]	LLMs	LLMs		
LLMGenPlan	[51]		LLMs		

## 2399 O.2 LTL Agent Interface

2400 Initially proposed as a formal verification for computer programs [40], Linear Temporal Logic (LTL) expressions  
 2401 have since then been extensively used in robotics [35–37] as a formalism that enables the extraction of guarantees  
 2402 on robot performance given a robot model, a high-level description of robotic actions, and a class of admissible  
 2403 environments. Recent work in embodied agent planning has adopted LTL expressions for their agent interface  
 2404 due to their high expressiveness and formality, bridging the gap between natural language and robotic control  
 2405 language [38, 39, 15]. Among these works [38] and [39] directly translate natural language actions and task goals  
 2406 into LTL to perform planning tasks, while [15] utilizes an intermediary representation called Signal Temporal  
 2407 Logic (STL) derived from LTL for action and subgoal planning. Due to the compactness and expressiveness of  
 2408 LTL, our work leverages LTL as a unified interface between modules and agents to describe task specifications.

2409 **P Annotation Details**

2410 **P.1 BEHAVIOR**

2411 For BEHAVIOR-100, we provide 2 additional sets of manual annotations: ground truth action sequences and  
2412 natural language task descriptions. The annotation is done based on observing the real demonstrations [1].

2413 **P.1.1 Annotating Action Sequence based on VR Human Demonstrations**

2414 The BEHAVIOR dataset consists of 100 tasks defined using the Behavior Domain Definition Language (BDDL),  
2415 which is similar to the PDDL format. Each BDDL file has three parts: a list of objects, initial states, and goal  
2416 states. To build the transition model for our simulator, we annotated the ground-truth action sequence for each  
2417 task. The annotation process involves an automatic pipeline for objects' states tracking and segmentation from  
2418 the demo, followed by manual effort for mapping state changes to action sequences.

2419 **P.1.2 Complicated Goals with Quantifiers**

2420 BDDL goals may contain quantifiers, such as `forpairs( ?jar.n.01 - jar.n.01)`, which need to be  
2421 translated into grounded goals, e.g., `((inside apple.n.01_1 jar.n.01_1) (inside`  
2422 `apple.n.01_2 jar.n.01_2))`. There can be different grounded goals that satisfy the same BDDL goal,  
2423 which we refer to as goal options. For example, `((inside apple.n.01_2 jar.n.01_1) (inside`  
2424 `apple.n.01_1 jar.n.01_2))` also satisfies the aforementioned BDDL goal. In general, one BDDL goal  
2425 may have a number of goal options, and each goal option has a number of grounded atomic goals.

2427 **P.1.3 Annotating Transition Models**

2428 **PDDL Problem File Annotation.** We devised an automatic pipeline to generate PDDL problem files from  
2429 BDDL files defined for tasks in BEHAVIOR . Specifically, we would sample a solvable goal option from the  
2430 BDDL goals, and select max \_num grounded goals from this option as the goal for the PDDL problem file,  
2431 then collect relevant objects and initial conditions for the grounded goals. This is done to ensure that the PDDL  
2432 planner could solve each problem within the desired time frame.

2433 **PDDL Domain File Annotation.** The PDDL domain file is modified from an existing PDDL domain file  
2434 written for BEHAVIOR <sup>¶</sup>. To make it compatible with our problem file and run in a PDDL planner, we modified  
2435 the naming for object types, added new predicates, and changed the preconditions and post-effects for some  
2436 operators.

2437 **P.2 VirtualHome**

2438 We build our annotation of VirtualHome on top of RobotHow [2], which provides a categorized list of household  
2439 tasks, each with a varying number of instructions. Each instruction includes a natural language goal description,  
2440 a VirtualHome action script, and the initial and final states of the scene. However, RobotHow does not offer  
2441 precise symbolic goals for each task, which limits the accuracy of evaluation due to noisy final states. For  
2442 instance, it is more precise to denote the final goal of the task “turn on light” as `on(light)`, rather than a noisier  
2443 version like `ontop(agent, chair) and on(light)`.

2444 To enhance the reliability and standardization of evaluation outcomes, we manually annotated 338 task instruc-  
2445 tions across 30 task categories. Our annotation process involved two main steps:

2446 **P.2.1 Wildcard Representations: Deriving Common Task-Related Final States**

2447 We identified common final states for all instructions within the same task category by using wildcards and  
2448 object properties. This allows for multiple solutions of each task. For example, we use the property `clothes` to  
2449 represent both `pants` and `shirts`. This step provides each task with a wildcard representation.

2450 **P.2.2 Extracting and Annotating Abstract Goals**

2451 Based on the wildcard representations, we extracted all concrete goals and manually added any missing ones  
2452 (which are usually few). We annotated the following three types of goals:

- 2453 • State Goals: Represent object and agent states in the current scene, such as `plugged_in(TV)`.
- 2454 • Relation Goals: Represent the spatial relationships that should be satisfied at the end of execution,  
2455 such as `ontop(agent, chair)`.

<sup>¶</sup>[https://github.com/wmccclinton/downward/blob/main/behavior\\_full.pddl](https://github.com/wmccclinton/downward/blob/main/behavior_full.pddl)

- 2456 • Action Goals: Represent actions that must be performed as part of the task instruction, especially those  
 2457 without post-effects. For example, in the task “*pat cat*”, the agent is required to perform the action  
 2458 *touch* even though it has no post-effect in the final scene state.

### 2459 **P.2.3 Grounding to Concrete Goals from Abstract Goals**

2460 In the VirtualHome simulator, each goal consists of an abstract goal and a concrete goal. An abstract goal is  
 2461 designed to be general and applicable to all potential trajectories. It may contain wildcards to represent similar  
 2462 objects, such as *book*, *notebook*, or *novel*. On the other hand, concrete goals are grounded to specific objects  
 2463 with their unique IDs, making them concrete and actionable.

2464 The conversion between an abstract goal and a concrete goal involves mapping the wildcards to relevant objects  
 2465 in the scene. This process transforms the abstract goals into specific, tangible objectives that can be directly  
 2466 acted upon by the embodied agent. By replacing the wildcards with the appropriate object IDs, the abstract goals  
 2467 are instantiated into concrete goals that are tailored to the specific environment and the available objects. For  
 2468 example, an abstract goal like “pick up a *book*” may be converted into a concrete goal such as “pick up *book\_1*”  
 2469 or “pick up *novel\_2*”, depending on the specific objects present in the scene and their corresponding IDs. This  
 2470 conversion process allows the embodied agent to have a clear understanding of the exact objects it needs to  
 2471 interact with to achieve the desired goal.

2472 The mapping of wildcards to relevant objects is a crucial step in bridging the gap between the high-level, abstract  
 2473 goals and the low-level, executable actions. It enables the embodied agent to ground the goals in the context of  
 2474 the specific environment it operates in, making the goals more precise and achievable.

### 2475 **P.2.4 Annotating Transition Models**

2476 In addition to the goals, we also annotate the transition model of preconditions and post-effects for each action  
 2477 in VirtualHome using standard PDDL formulations. Since VirtualHome does not provide an existing PDDL  
 2478 domain file, we define the predicate list and provide a PDDL implementation for each action.

2479 The predicate list includes all relations (e.g., *inside*, *facing*), object states (e.g., *plugged\_in*, *closed*), and object  
 2480 properties related to actions (e.g., *grabbable*, *has\_switch*). The PDDL version of the actions follows the same  
 2481 preconditions and effects as their implementation in VirtualHome Github. The logic involves *and*, *or*, *not*, *when*,  
 2482 *exists*, and *forall*. Approximately 10 actions in VirtualHome simulate human body movement and have no effects  
 2483 on the environment. For example, the action *read* requires the robot to hold something readable but has no  
 2484 effects. Such actions would possibly be left out in the success rate by planner metrics because they would not  
 2485 contribute to the achievement of goals. However, we still provide their PDDL definitions and task the LLM to  
 2486 predict them, which is evaluated by the logic form accuracy metric.

## 2487 **Q Simulator Implementation Details**

### 2488 **Q.1 BEHAVIOR Implementation Details**

2489 We developed an environment for evaluating LLMs in behavioral tasks, which includes two types of simulators.  
 2490 The first simulator is based on the iGibson framework [75], supporting visual rendering and physical simulations  
 2491 via PyBullet. The second simulator is symbolic, recording activities in a graph-based format. Both simulators  
 2492 share the same set of actions and object states.

#### 2493 **Q.1.1 Building the Symbolic Simulator and Transition Model Implementation**

2494 Details about the arguments, preconditions, and post-effects of each action are presented in Table 24. An object’s  
 2495 interactability is defined as: (1) The object is not enclosed within another object, and (2) The object is within the  
 2496 agent’s reach. However, as we currently allow auto-navigation for the agent, which automatically invokes the  
 2497 NAVIGATE\_TO action, this reachability requirement is not necessary for LLMs to fulfill.

Table 24: BEHAVIOR Symbolic Simulator Implementation: Action Transition Models.

Action Name	Arguments	Preconditions	Post Effects
NAVIGATE_TO	tar_obj1	tar_obj1 is interactable.	The agent is next to tar_obj1.
LEFT_GRASP	tar_obj1	tar_obj1 is interactable; One of the agent’s hands is empty; tar_obj1 is small enough.	tar_obj1 is in the agent’s left hand.

Continued on next page

Table 24 – continued from previous page

Action Name	Arguments	Precondition	Post Effects
RIGHT_GRASP	tar_obj1	tar_obj1 is interactable; The agent's right hand is empty; tar_obj1 is small enough.	tar_obj1 is in the agent's right hand.
LEFT_RELEASE	tar_obj1	tar_obj1 is interactable; The agent's left hand holds tar_obj1.	tar_obj1 is released to the floor.
RIGHT_RELEASE	tar_obj1	tar_obj1 is interactable; The agent's right hand holds tar_obj1.	tar_obj1 is released to the floor.
LEFT_PLACE_ONTOP	tar_obj1	tar_obj1 is interactable; The agent's left hand holds an object.	The object in the agent's left hand is on top of tar_obj1.
RIGHT_PLACE_ONTOP	tar_obj1	tar_obj1 is interactable; The agent's right hand holds an object.	The object in the agent's right hand is on top of tar_obj1.
LEFT_PLACE_INSIDE	tar_obj1	tar_obj1 is interactable; The agent's left hand holds an object; tar_obj1 is big enough and open if openable.	The object in the agent's left hand is inside tar_obj1.
RIGHT_PLACE_INSIDE	tar_obj1	tar_obj1 is interactable; The agent's right hand holds an object; tar_obj1 is big enough and open if openable.	The object in the agent's right hand is inside tar_obj1.
LEFT_PLACE_NEXTTO	tar_obj1	tar_obj1 is interactable; The agent's left hand holds an object.	The object in the agent's left hand is next to tar_obj1.
RIGHT_PLACE_NEXTTO	tar_obj1	tar_obj1 is interactable; The agent's right hand holds an object.	The object in the agent's right hand is next to tar_obj1.
LEFT_PLACE_UNDER	tar_obj1	tar_obj1 is interactable; The agent's left hand holds an object.	The object in the agent's left hand is under tar_obj1.
RIGHT_PLACE_UNDER	tar_obj1	tar_obj1 is interactable; The agent's right hand holds an object.	The object in the agent's right hand is under tar_obj1.
LEFT_TRANSFER_CONTENTS_INSIDE	tar_obj1	tar_obj1 is interactable; The agent's left hand holds an object; tar_obj1 is big enough and open if openable.	The contents inside the object in the agent's left hand are in tar_obj1.
RIGHT_TRANSFER_CONTENTS_INSIDE	tar_obj1	tar_obj1 is interactable; The agent's right hand holds an object; tar_obj1 is big enough and open if openable.	The contents inside the object in the agent's right hand are in tar_obj1.
LEFT_TRANSFER_CONTENTS_ONTOP	tar_obj1	tar_obj1 is interactable; The agent's left hand holds an object.	The contents inside the object in the agent's left hand are on top of tar_obj1.
RIGHT_TRANSFER_CONTENTS_ONTOP	tar_obj1	tar_obj1 is interactable; The agent's right hand holds an object.	The contents inside the object in the agent's right hand are on top of tar_obj1.
LEFT_PLACE_NEXTTO_ONTOP	tar_obj1, tar_obj2	tar_obj1 and tar_obj2 are interactable; The agent's left hand holds an object.	The object in the agent's left hand is next to tar_obj1 and on top of tar_obj2.

Continued on next page

Table 24 – continued from previous page

Action Name	Arguments	Precondition	Post Effects
RIGHT_PLACE_NEXTTO_ONTOP	tar_obj1, tar_obj2	tar_obj1 and tar_obj2 are interactable; The agent's right hand holds an object.	The object in the agent's right hand is next to tar_obj1 and on top of tar_obj2.
TOGGLE_ON	tar_obj1	tar_obj1 is interactable; One of the agent's hands is empty; tar_obj1 is toggled off and closed.	tar_obj1 is toggled on.
TOGGLE_OFF	tar_obj1	tar_obj1 is interactable; One of the agent's hands is empty; tar_obj1 is toggled on.	tar_obj1 is toggled off.
CLOSE	tar_obj1	tar_obj1 is interactable; One of the agent's hands is empty; tar_obj1 is open.	tar_obj1 is closed.
OPEN	tar_obj1	tar_obj1 is interactable; One of the agent's hands is empty; tar_obj1 is closed and not toggled on.	tar_obj1 is open.
CLEAN	tar_obj1	tar_obj1 is interactable; tar_obj1 is dusty or stained; the agent has a cleaning tool.	tar_obj1 is not dusty or stained (according to the tool).
DRY	tar_obj1	tar_obj1 is interactable; tar_obj1 is soaked.	tar_obj1 is dry.
SLICE	tar_obj1	tar_obj1 is interactable; tar_obj1 is not sliced; the agent has a knife.	tar_obj1 is sliced.
SOAK	tar_obj1	tar_obj1 is interactable; One of the agent's hands is empty; tar_obj1 is dry and in a toggled on sink or a pot.	tar_obj1 is soaked.
FREEZE	tar_obj1	tar_obj1 is interactable; One of the agent's hands is empty; tar_obj1 is unfrozen and in the fridge.	tar_obj1 is frozen.
UNFREEZE	tar_obj1	tar_obj1 is interactable; tar_obj1 is frozen.	tar_obj1 is unfrozen.
COOK	tar_obj1	tar_obj1 is interactable; One of the agent's hands is empty; tar_obj1 is not cooked; tar_obj1 is on/in the pan.	tar_obj1 is cooked.

2498 **Q.1.2 BEHAVIOR State Space and Action Space**

2499 Our environment defines 15 object states, derived from a subset of iGibson's state definitions, and offers 30  
 2500 actions that agents can use to interact with the objects in a task scene. The names of these actions and object  
 2501 states are provided in Table 25. Object states are categorized into unary states, describing conditions or attributes  
 2502 of an object, and binary states, representing physical relationships between two objects. Actions are divided  
 2503 into state actions, which modify the unary states of objects, and spatial actions, which modify the binary states  
 2504 between two objects or between an object and the agent.

2505 **Q.2 VirtualHome Implementation Details**

2506 **Q.2.1 VirtualHome State Space and Action Space**

2507 We have developed an evaluation environment for LLMs using the VirtualHome simulator as our foundation. At  
 2508 the core of this environment is a MotionPlanner, which leverages the EnvironmentGraph from VirtualHome  
 2509 to record the state of the environment. Additionally, the MotionPlanner incorporates runtime error analysis  
 2510 for all 42 possible actions. To ensure accurate execution outcomes, we made necessary modifications to the

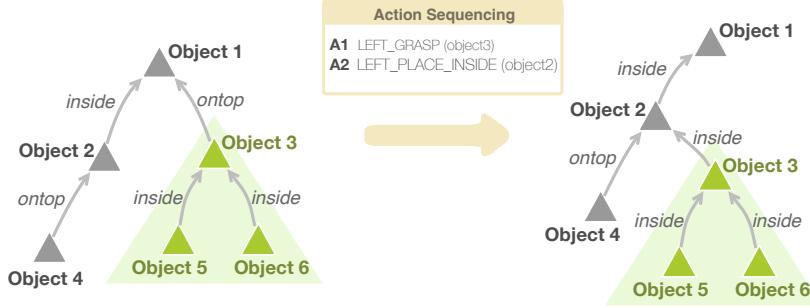


Figure 58: Implementation of Kinetics Tree for the Transition Model in BEHAVIOR .

Table 25: BEHAVIOR State Space and Action Space.

State Space		Action Space	
Unary States	Binary States	State Actions	Spatial Actions
<i>Cooked</i>	<i>Inside</i>	TOGGLE_ON	LEFT_GRASP
<i>Dusty</i>	<i>OnFloor</i>	TOGGLE_OFF	RIGHT_GRASP
<i>Frozen</i>	<i>OnTop</i>	CLOSE	LEFT_RELEASE
<i>Open</i>	<i>Under</i>	OPEN	RIGHT_RELEASE
<i>Sliced</i>	<i>NextTo</i>	CLEAN	LEFT_PLACE_ONTOP
<i>Soaked</i>		DRY	RIGHT_PLACE_ONTOP
<i>Stained</i>		SLICE	LEFT_PLACE_NEXTTO
<i>ToggledOn</i>		SOAK	RIGHT_PLACE_NEXTTO
<i>Slicer</i>		FREEZE	LEFT_PLACE_INSIDE
<i>CleaningTool</i>		UNFREEZE	RIGHT_PLACE_INSIDE
		COOK	LEFT_PLACE_UNDER RIGHT_PLACE_UNDER LEFT_TRANSFER_CONTENTS_INSIDE RIGHT_TRANSFER_CONTENTS_INSIDE LEFT_TRANSFER_CONTENTS_ONTOP RIGHT_TRANSFER_CONTENTS_ONTOP LEFT_PLACE_NEXTTO_ONTOP RIGHT_PLACE_NEXTTO_ONTOP NAVIGATE_TO

2511 VirtualHome simulator, aligning its action executors with expected results. The current state space and action  
 2512 space utilized in this environment are detailed in Table 26.

### 2513 Q.3 Evaluation Settings of LLMs

#### 2514 Q.3.1 Decoding Parameters

2515 To ensure standardization and consistency across all models, we utilized the same set of decoding parameters  
 2516 for all the LLMs evaluated in this study. Specifically, we used a temperature of zero for all models, as our goal  
 2517 was to use the arg max under the model’s distribution. Furthermore, since several of the models only support  
 2518 temperature-based sampling and not other sampling methods, we limited ourselves to temperature-scaling during  
 2519 the sample process. It is important to note that for a given prompt, the model’s completion involves sampling,  
 2520 which introduces randomness in determining the specific completion decoded for each instance. However, in our  
 2521 scenarios, this is not a significant factor since we are decoding the arg max through low-temperature decoding.

#### 2522 Q.3.2 Evaluation Cost

2523 To perform a single run of all models on our benchmark, a total of 180 runs would be required. This involves  
 2524 evaluating each specific model on each specific simulator ability module. The total cost of this process amounts  
 2525 to approximately 126,900,000 tokens and 50,760 queries across all models, which sums to a compute cost of  
 2526 \$679.24. For open-source models, the costs are calculated based on the pricing of the Together AI API <sup>1</sup>.

<sup>1</sup><https://www.together.ai/>

Table 26: VirtualHome State Space and Action Space.

State Space		Action Space	
Unary States	Binary States	State Actions	Spatial Actions
<i>Closed</i>	<i>On</i>	OPEN	TURN_TO
<i>Open</i>	<i>Inside</i>	SIT	PUT_BACK
<i>On</i>	<i>Between</i>	STANDUP	PUT_IN
<i>Off</i>	<i>Close</i>	SLEEP	POUR
<i>Sitting</i>	<i>Facing</i>	WAKEUP	PUT_ON
<i>Dirty</i>	<i>Holds_RH</i>	CLOSE	FIND
<i>Clean</i>	<i>Holds_LH</i>	DRINK	RUN
<i>Lying</i>		GRAB	WALK
<i>Plugged_in</i>		LOOKAT	POINT_AT
<i>Plugged_out</i>		LOOKAT_SHORT	TOUCH
		LOOKAT_LONG	WATCH
		SWITCH_OFF	MOVE
		SWITCH_ON	RELEASE
		TYPE	DROP
		PUSH	
		PULL	
		SQUEEZE	
		WASH	
		RINSE	
		SCRUB	
		EAT	
		PLUG_IN	
		PLUG_OUT	
		CUT	
		READ	
		LIE	

2527 **Q.3.3 Models**

2528 Due to space limitation, we use shorthand model names in the main paper. Here we provide the details of the  
 2529 models in Table 9.

2530