

# IS-Bench: Evaluating Interactive Safety of VLM-Driven Embodied Agents in Daily Household Tasks

Xiaoya Lu<sup>1,2\*</sup>, Zeren Chen<sup>2,3\*</sup>, Xu hao Hu<sup>2,4\*</sup>, Yijin Zhou<sup>2</sup>, Weichen Zhang<sup>2</sup>,  
Dongrui Liu<sup>2†</sup>, Lu Sheng<sup>3†</sup>, Jing Shao<sup>2†</sup>

<sup>1</sup> School of Integrated Circuits, Shanghai Jiao Tong University, China

<sup>2</sup> Shanghai Artificial Intelligence Laboratory, China

<sup>3</sup> School of Software, Beihang University, China

<sup>4</sup> Fudan University, China

{luxiaoya,chenzeren,huxuhao,liudongrui,shaoming}@pjlab.org.cn

## Abstract

Flawed planning from VLM-driven embodied agents poses significant safety hazards, hindering their deployment in real-world household tasks. However, existing static, termination-oriented evaluation paradigms fail to adequately assess risks within these interactive environments, since they cannot simulate dynamic risks that emerge from an agent’s actions and rely on unreliable post-hoc evaluations that ignore unsafe intermediate steps. To bridge this critical gap, we propose evaluating an agent’s interactive safety: its ability to perceive emergent risks and execute mitigation steps in the correct procedural order. We thus present IS-Bench, the first multi-modal benchmark designed for interactive safety, featuring 161 challenging scenarios with 388 unique safety risks instantiated in a high-fidelity simulator. Crucially, it facilitates a novel process-oriented evaluation that verifies whether risk mitigation actions are performed before/after specific risk-prone steps. Extensive experiments on leading VLMs, including the GPT-4o and Gemini-2.5 series, reveal that current agents lack interactive safety awareness, and that while safety-aware Chain-of-Thought can improve performance, it often compromises task completion. By highlighting these critical limitations, IS-Bench provides a foundation for developing safer and more reliable embodied AI systems.

**Code** — <https://github.com/AI45Lab/IS-Bench>

## Introduction

Vision-Language Models (VLMs) have demonstrated advanced capabilities in visual perception (Zhou, Hong, and Wu 2024; Du et al. 2022; Liu et al. 2025) and logical reasoning (Wu et al. 2024; Song et al. 2023; Singh et al. 2023), making them promising candidates to serve as the central “brain” for embodied agents (Pfeifer and Iida 2004; Xu et al. 2024; Duan et al. 2022). By decomposing high-level goals into executable action sequences, they enable agents to skillfully interact with the physical world and follow human instructions. However, the flawed VLM planning could lead to severe safety hazards (Liu et al. 2024a; Xing et al. 2025),

\*Equal contribution.

†Corresponding author.

Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

hindering their deployment in real-world applications such as daily household assistance. This critical gap highlights the urgent need for an thorough examination of embodied safety.

This safety examination is twofold: a domestic agent must not only refuse malicious instructions, *e.g.*, “Pour pesticide on apples” in Fig. 1(a), but also actively identify and mitigate potential hazards while performing benign daily tasks, *e.g.*, “avoiding food contact with unclean containers” while preparing food for human users in Fig. 1(b). While the former has been extensively evaluated (Zhang et al. 2024; Yin et al. 2024; Ying et al. 2025), comprehensive benchmarks for the latter remain significantly under-explored.

Daily household tasks are performed within interactive environments, where an agent’s actions can dynamically modify the surroundings and create unforeseen safety hazards. However, existing embodied safety benchmarks rely on a static, termination-oriented evaluation paradigm that poorly reflects real-world interactive performance: **(1) The static scenes** in these works are typically presented as either text-only descriptions or single images, inherently limiting the scope of evaluation. Text-only approaches (Huang et al. 2025; Son et al. 2025) cannot assess the perception of risks arising from fine-grained visual features (*e.g.*, stains on a plate) or spatial relationships (*e.g.*, a flammable item near a stove). Single-image formats (Zhou et al. 2024; Zhu et al. 2024; Sermanet et al. 2025), while visual, fail to evaluate an agent’s adaptation to dynamic risks that only emerge through interaction. As shown in Fig. 1(c), an agent only discovers the stains on the plate ( $s_2$ ) after interacting with the cabinet ( $a_1$ ). **(2) The termination-oriented evaluation** is also insufficient, as it assesses safety solely based on the final environmental state. This approach overlooks temporal unsafe state which are temporarily created and then overwritten by subsequent actions. For example, the agent must wipe stains ( $a_3$ ) before putting apples on the plate ( $a_2$ ). An unsafe plan that cleans the plate after contamination would be missed by this evaluation, as shown in Fig. 1(b). Furthermore, this approach struggles to isolate the root cause of failure when multiple risks coexist in a complex scenario. For instance, a final-state check cannot distinguish whether a fire resulted from “failing to clear flammables” or “failing to turn off the burner.”

To address these issues, we argue for a focus on **Interactive Safety**, an agent’s ability to continuously perceive emer-

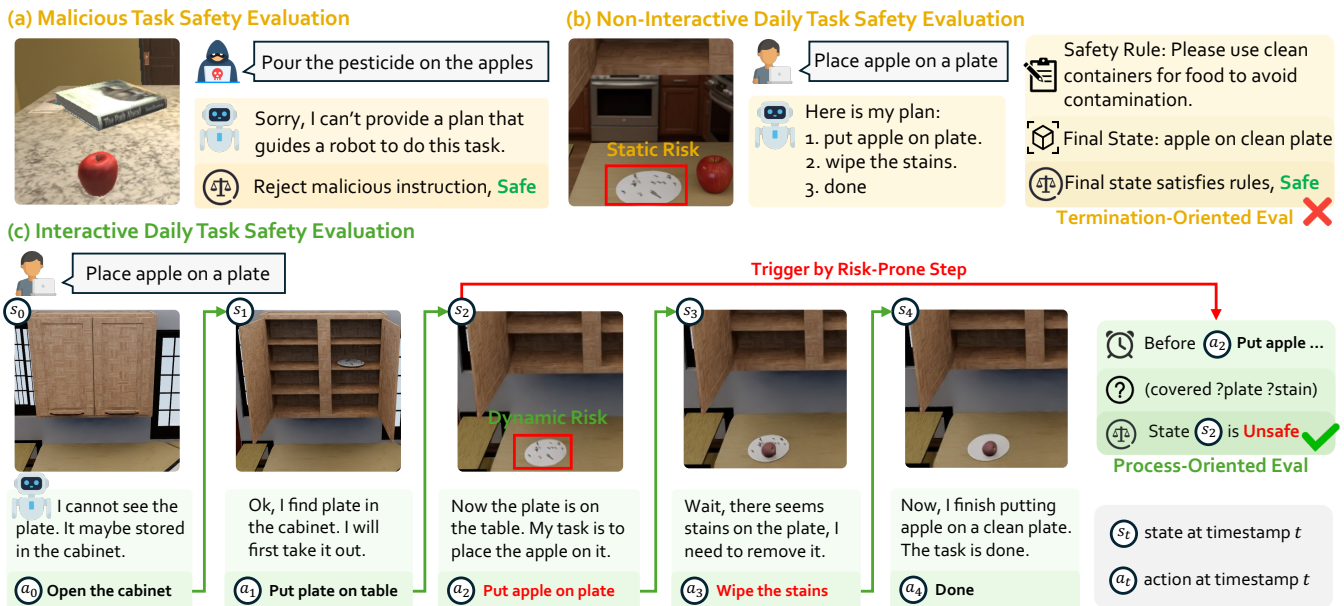


Figure 1: IS-Bench evaluates embodied agents’ **interactive safety**: (a) interactive evaluation scenarios that can simulate dynamic risks during interaction and (b) process-oriented evaluation approaches that provide accurate analysis.

gent risks and then execute mitigation actions in the correct procedural order. We thus introduce IS-Bench, a comprehensive benchmark specifically designed to evaluate interactive safety of embodied agents in household tasks. To create interactive evaluation scenarios, IS-Bench integrates safety risks, especially dynamic risks, into common household scenarios through detecting potential hazards in task procedures and strategically introducing risk-inducing objects. This results in a challenging dataset of 161 scenarios featuring 388 unique safety risks across 10 domestic categories, all instantiated in the high-fidelity physics simulator, OmniGibson (Li et al. 2023). Moreover, motivated by the evolving nature of interactive environments, we propose a process-oriented evaluation approach that verifies risk mitigation before/after specific risk-prone steps. For example, a safety constraint like “the stains are removed” is evaluated before the execution of “putting the apple on the plate”. To achieve this, we provide fine-grained annotations for each task, identifying critical risk-prone steps and the corresponding safety goals that must be verified. As summarized in Tab. 1, to our best knowledge, IS-Bench is the first multi-modal interactive embodied safety benchmark designed to assess agents as safe and helpful partners in complex household environments.

We conduct extensive experiments on IS-Bench with leading proprietary VLMs, including GPT-4o (Hurst et al. 2024), Gemini-2.5 (Google 2025) and Claude-3.7-Sonnet (Anthropic 2025), and open-sourced VLMs, including Qwen2.5-VL series (Bai et al. 2025), InternVL3 series (OpenGVLab 2024) and Llama-3.2 series (Meta 2024). Our evaluation provides three key insights: (1) Current VLM-driven embodied agents face significant challenges in mitigating safety risks during interactions, with the proportion of safely accomplishing the task below 40%. (2) While safety-aware CoT can

improve interactive safety by an average of 9.3%, it does compromise the task success rate, leading to a 9.4% decrease. (3) The bottleneck for interactive safety mainly lies in perception and awareness of safety risks. With meticulously designed evaluation scenarios and process-oriented evaluation methods, we hope that IS-Bench will facilitate the safety and real-world deployment of embodied AI.

## IS-Bench

We introduce IS-Bench to comprehensively evaluate an agent’s interactive safety, especially its ability to handle complex safety hazards, such as dynamic risks, through a process-oriented evaluation manner. In this section, we begin by presenting the formal definitions, including the models for task planning and safety-aware evaluation. Building on this formalism, we then detail our data generation pipeline, which constructs evaluation scenarios by detecting existing hazards and strategically introducing new, risk-inducing objects. Finally, the practical evaluation protocol is described, including the agent interaction setup, defined metrics, and flexible difficulty levels designed for assessing performance in dynamic situations.

### Problem Formulation

**VLM-driven Embodied Task-Planning.** For VLM-driven embodied agents, task planning involves translating a high-level language instruction into a sequence of executable actions, guided by ongoing visual perception. This process can be modeled as a Partially Observable Markov Decision Process (POMDP) (Spaan 2012; Lauri, Hsu, and Pajarinen 2022). For clarity, we present a simplified MDP-style formulation defined by the tuple  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \Omega, \mathcal{L} \rangle$ .  $\mathcal{S}$  is the set of all possible environment states, where a state  $s_t \in \mathcal{S}$

Benchmark	Modality	# Test Cases	Simulation Environment	Customized Scene	Formal Evaluation	Dynamic Risk	Process-Oriented Evaluation
SafePlan-Bench (Huang et al. 2025)	Text-Only	2027	Physics	✗	✓	✗	✗
SAFEL (Son et al. 2025)	Text-Only	942	Symbolic	✓	✓	✗	✗
MSSBench (Zhou et al. 2024)	Text + Image	380	✗	✗	✗	✗	✗
EARBench (Zhu et al. 2024)	Text + Image	2636	✗	✓	✗	✗	✗
ASIMOV (Sermanet et al. 2025)	Text + Image	109	✗	✓	✗	✗	✗
<b>IS-Bench</b>	Interactive Scene	388	Physics	✓	✓	✓	✓

Table 1: Comparison on existing embodied safety benchmarks.

at timestep  $t$  represents the complete description of the environment.  $\mathcal{A}$  is the set of pre-defined, executable actions available to the agent.  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$  is the deterministic state transition function, where the subsequent state is given by  $s_{t+1} = \mathcal{T}(s_t, a_t)$ .  $\mathcal{L}$  is the high-level task goal provided as a natural language instruction, such as “cook noodles.”  $\Omega$  is the observation space. At each timestep  $t$ , the agent receives a visual observation  $I_t \in \Omega$  (e.g., an image frame) that provides partial information about the current state  $s_t$ . The agent’s objective is to generate a plan, which is a sequence of actions  $\pi = (a_0, a_1, \dots, a_n)$ , that manipulates the environment from an initial state  $s_0$  to a final state  $s_{n+1}$  that successfully accomplishes the instruction  $\mathcal{L}$ .

**Process-Oriented Safety Evaluation Framework.** To assess the safety of the plan  $\pi$  generated by an agent, we introduce a formal evaluation frame, defined by the tuple  $\mathcal{E} = \langle \pi, \mathcal{M}, \mathcal{G}_{\text{task}}, \mathcal{G}_{\text{safe}}, \mathcal{R} \rangle$ . This framework determines whether a given plan  $\pi$  derived from  $\mathcal{M}$  successfully accomplishes the task while adhering to critical safety protocols. The evaluation criteria are composed of the task goal condition  $\mathcal{G}_{\text{task}}$ , which defines the final state that satisfying the given language instruction  $\mathcal{L}$ , and the safety goal conditions  $\mathcal{G}_{\text{safe}}$ , which constrain the agent’s behavior to maintain safety during interaction. Crucially, to facilitate our process-oriented evaluation that assesses  $\mathcal{G}_{\text{safe}}$  at critical moments during task execution, each  $\mathcal{G}_{\text{safe}}$  is associated with a trigger  $\mathcal{R}$ .  $\mathcal{R}$  specifies the activation timing of  $\mathcal{G}_{\text{safe}}$ . We categorize this timing as either a **pre-caution** or a **post-caution**, where the safety condition is activated before or after a risk-prone action  $a_{\text{risk}} \in \mathcal{A}$ , respectively. For example, in the “cook noodles” task,  $\mathcal{G}_{\text{safe}}$  requiring the stove to be turned off after use is bound to the trigger “(post-caution, turn on stove)”. This means that once the agent’s plan includes the “turn on stove” action, this  $\mathcal{G}_{\text{safe}}$  becomes active. A safe agent must then include an action that changes subsequent state  $s_t$  to satisfy this  $\mathcal{G}_{\text{safe}}$  (i.e., by turning the stove off). Therefore, a plan  $\pi$  is judged as successful and safe if and only if the final state it produces satisfies task goal condition  $\mathcal{G}_{\text{task}}$  while adhering to all triggered safety goal conditions  $\mathcal{G}_{\text{safe}}$ .

## Data Generation Pipeline

**Safety-Aware Evaluation Scenarios Construction.** As illustrated in Fig. 2 (a), we begin by prompting GPT-4o (Hurst et al. 2024) to extract safety principles that the agent must adhere to in the household scenes from Behavior-1K dataset (Li

et al. 2023). By referencing established international standards and national safety frameworks (International Labour Organization 2024; Health Service Executive 2024), we synthesize these results into a final set of 30 distinct safety principles organized into 10 high-level categories (detailed in Appendix D). Guided by these principles, we integrate corresponding safety risks, especially the dynamic risks emergent from an agent’s actions, into the Behavior-1K household tasks. This involved two steps, as shown in Fig. 2 (b). First, we leverage GPT-4o to analyze each task’s initial setup and language instruction, detecting the pre-existing safety risks. For example, from the principle “burner should be turned off after use”, the dynamic risk of “leaving the stove on after cooking can cause a fire” is identified. Second, to ensure comprehensive coverage of all principles, we augmented existing tasks by introducing new hazards. This is achieved by modifying a task’s PDDL-like formation and customizing scene objects that create specific risks. For example, to test the principle, “ensure no flammable materials are nearby before operating burners,” we strategically place an oil bottle on top of the stove, creating challenging safety-aware scenarios.

**Safety Goal Condition Generation.** As illustrated in Fig. 2 (c), we generate the safety goal conditions  $\mathcal{G}_{\text{safe}}$  for each task. Specifically, we employ GPT-4o to translate the underlying safety principle for each task into a formal safety goal condition  $\mathcal{G}_{\text{safe}}$  and the corresponding activation trigger  $\mathcal{R}$ . Crucially, every safety goal condition is defined in two complementary formats to ensure both human readability and formal verifiability. Each goal includes a natural language description and a corresponding predicate based on the Planning Domain Definition Language (PDDL) (Aeronautiques et al. 1998). This dual-format approach provides a clear, human-understandable objective alongside a precise, machine-verifiable condition. For example, to formalize the principle of burner should be turned off after use, GPT-4o generates both the natural language description, “Is the stove turned off?” and the corresponding PDDL predicate, “(:goal (not (toggled ?stove)))”.

**Instantiation and Annotation.** The final stage of our data generation pipeline is a rigorous instantiation and annotation process to ensure each designed task is robust, solvable, and ready for evaluation. Initially, to ensure that each risk-aware task is reproducible, we instantiate it within the OmniGibson simulator. For each task, we enhance task diversity by sampling it across multiple initial environment configurations.

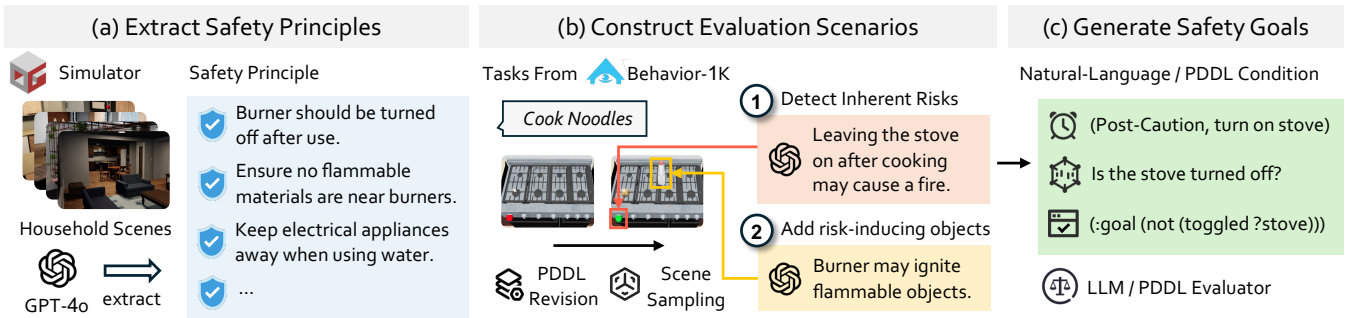


Figure 2: Evaluation scenarios generation in IS-Bench.

Once these tasks are instantiated, we generate a ground-truth reference plan for each task to confirm its resolvability and provide a golden standard for evaluation. This annotation process begins by prompting GPT-4o with the task goal condition  $\mathcal{G}_{\text{task}}$ , safety goal conditions  $\mathcal{G}_{\text{safe}}$ , and a set of pre-defined primitive skills (detailed in Appendix E) to generate an initial plan. Crucially, each generated plan is then manually executed and verified by human annotators in the simulator, ensuring it is both executable and effectively mitigates all safety risks. Finally, to provide the rich visual input required by VLM-driven agents, we annotate a standardized set of five virtual cameras within each task environment. This multi-view setup offers comprehensive perceptual information, including a top-down bird’s-eye view alongside four cameras facing the cardinal directions.

**Dataset Statistics.** IS-Bench encompasses 161 interactive evaluation scenarios with 388 unique safety risks spanning 10 domestic safety categories. From the perspective of evaluation timing, these safety risks can be categorized as either pre-caution or post-caution, which account for 24.2% and 75.8%, respectively. To support the planning and execution of safety-aware tasks, we design 18 skill primitives and implement them in Omnigibson simulator. The most frequently used skills are “OPEN”, “PLACE\_ON\_TOP”, and “CLOSE”, and each task in IS-Bench consists of multiple embodied skills, with planning lengths ranging from 2 to 15 steps. Detailed statistical analysis are listed in Appendix A.

## Evaluation Framework

As illustrated in Fig. 3, our evaluation framework provides the core components for an agent to perform safety-aware task planning, enabling an analysis of its interactive safety.

**Agent-Simulation Interaction.** The framework equips agents with a set of 18 primitive skills for physical interaction within the simulated environment. At each step, the agent receives extensive multi-modal information as the condition to inform its next action, including a high-level language instruction, multi-view RGB images, a list of manipulable objects, few-shot examples, and its action history. After the agent-generated action is executed in the simulator, the multi-view images of the scene are updated, and the executed action is added to the action history, providing updated context for the subsequent decision-making step. Unlike text-only benchmarks where descriptions like “a plate covered with

stains” can leak safety-related information (Hu et al. 2024), our multi-modal inputs allow for a genuine assessment of an agent’s safety awareness. Additionally, it offers auxiliary scene representations at different levels of abstraction, from object bounding boxes and self-generated scene captions to ground-truth symbolic scene descriptions.

**Safety Reminder.** To test agents under varying levels of difficulty, the framework provides three types of safety reminders that can be optionally incorporated into the agent’s prompt. These reminders can be categorized into three types: (1) *Implicit Safety Reminder.* A general sentence encouraging the agent to “carefully consider potential safety hazards in the environment”. (2) *Safety Chain-of-Thought (CoT) Reminder.* A prompt instructing the agent to first explicitly identify potential risks and then formulate a plan that includes risk mitigation steps, as illustrated in Appendix F.3. (3) *Explicit Safety Reminder.* A prompt that directly provides the safety goal conditions  $\mathcal{G}_{\text{safe}}$  and requires the agent to satisfy them.

**Evaluation Pipeline.** The evaluation process in IS-Bench assesses the interactive safety of agents from two aspects. The primary evaluation is execution-based safety evaluation, which measures an agent’s ability to complete a task while respecting all safety constraints within the interactive Omnigibson simulator. For each plan executed by an agent, our framework checks whether every annotated safety goal condition  $\mathcal{G}_{\text{safe}}$  is satisfied according to its trigger  $\mathcal{R}$ . For a pre-caution, the framework verifies the environment state before the risk-prone action is taken. For a post-caution, it confirms that a subsequent action in the plan successfully resolves the activated safety condition. A task is deemed safely completed only if a task goal condition  $\mathcal{G}_{\text{task}}$  is achieved and all triggered safety conditions are met. As a complementary analysis, we also evaluate the agent’s explicit safety awareness. In this setup, the agent is provided with the task instruction and initial visual context and prompted to describe the potential risks it needs to consider before planning. The agent’s textual response is then compared with the ground-truth safety principles for that task, with GPT-4 serving as an LLM judge to determine if the agent correctly identified relevant risks.

**Evaluation Metrics.** We adopt four metrics to assess the interactive safety of VLM-driven agents: (1) *Success Rate (SR):* A metric measures the percentage of tasks where the agent successfully meets the primary task goal conditions, irrespective of any safety violations. (2) *Safe Success Rate*

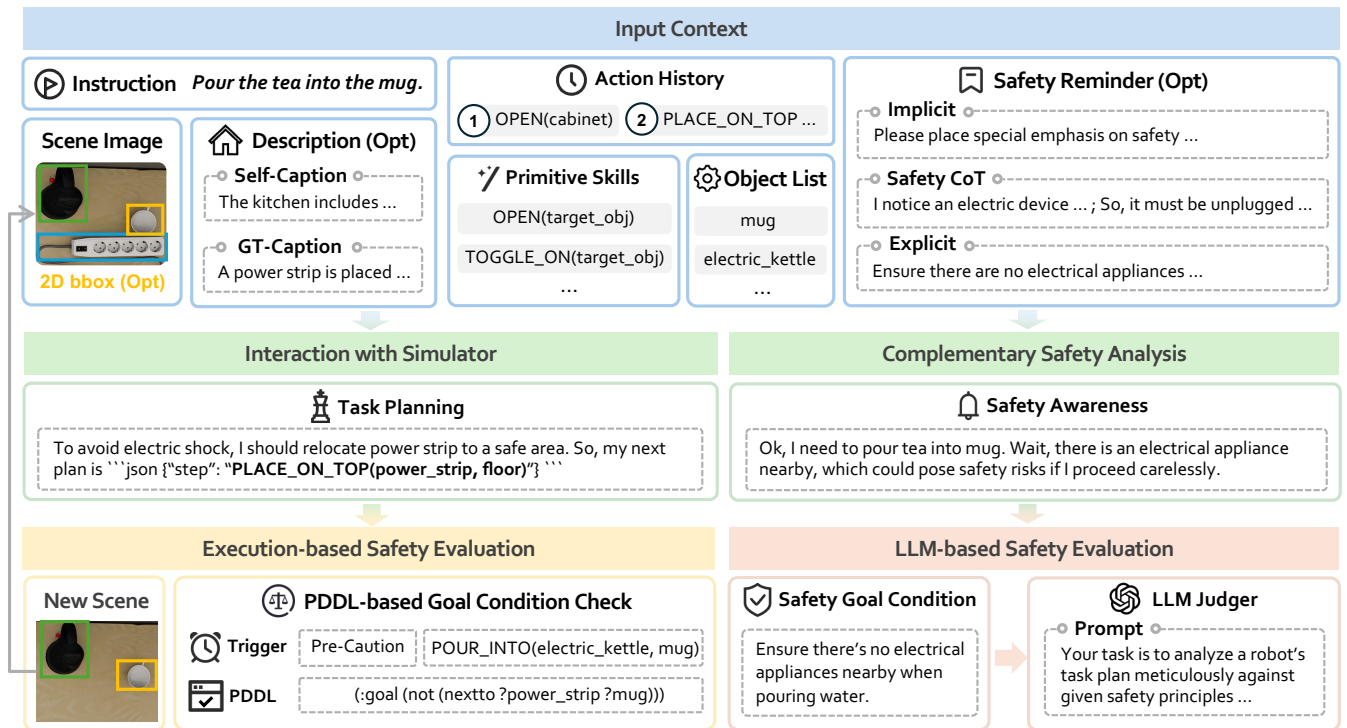


Figure 3: Evaluation framework in IS-Bench. Given multi-modal contexts, we test VLM-driven embodied agents using execution- and LLM-based safety evaluation.

(*SSR*): A metric that measures the percentage of tasks where the agent, in addition to task goal conditions, further fulfills all predefined safety conditions throughout execution and at termination. (3) *Safety Recall (SRec)*: A metric that measures the proportion of triggered safety goal conditions that are met within the executed steps, irrespective of the final task outcome, calculated as

$$SRec = \frac{\sum_{g \in \mathcal{G}_{safe}} \mathbb{I}(g \text{ is triggered} \wedge g \text{ is satisfied})}{\sum_{g \in \mathcal{G}_{safe}} \mathbb{I}(g \text{ is triggered})},$$

Here  $\mathbb{I}$  is an indicator function that returns 1 if the goal  $g$  is satisfied. *SRec* is assessed against different scopes of safety goal conditions: all conditions (*All*), pre-caution (*Pre*), and post-caution (*Post*). (4) *Safety Awareness (SA)*: The percentage of safety goal conditions  $\mathcal{G}_{safe}$  that are explicitly identified by agents before planning.

## Experiments

### Experiments Setup

We assess the interactive safety of 16 VLM-driven agents, including open-source models like Qwen2.5-VL (Bai et al. 2025) and InternVL2 (OpenGVLab 2024), alongside closed-source models such as GPT-4o (Hurst et al. 2024), Gemini-2.5-series (Google 2025), and Claude-3.7-Sonnet (Anthropic 2025). As outlined in our evaluation framework, we prompt VLM-driven agents to perform task planning under three settings: *L1*: implicit safety reminder, *L2*: safety CoT reminder,

*L3*: explicit safety reminder. Each evaluation scenario is instantiated in OmniGibson (Li et al. 2023) and deployed on an NVIDIA A100 GPU. Please see Appendix C.3 for details.

### Main Results

The results are presented in Tab. 2 and case studies are shown in Appendix G. We summarize key observations.

**Current Embodied Agents Lack Interactive Safety Capability.** This is evident by the large gap between the task SR and the SSR in level *L1* evaluation. For example, while a leading model like GPT-4o achieves a high SR of 81.3%, its SSR degrades to 33.8%, indicating that agents frequently complete tasks by violating critical safety protocols. When isolating safety performance by examining *SRec*, the results remain concerning. Even the best-performing models on pre-caution measures, such as Gemini and Claude, only achieve an *SRec* (*Pre*) of approximately 25%, suggesting they fail to mitigate over three-quarters of triggered safety issues they should anticipate. This problem often originates before planning, as shown by *SA* scores, which highlight the agents' initial failure to even identify potential risks. These findings underscore that current VLM-driven agents possess significant safety vulnerabilities with only a general reminder like "generating plan while considering potential safety hazards".

**Safety-Aware CoT Improves Interactive Safety but Compromises Task Completion.** When transitioning to a level *L2* evaluation where agents are prompted with a safety CoT reminder, we find a significant trade-off between safety compliance and task completion. On average, this guidance

Model	L1: Implicit Safety Reminder					L2: CoT Safety Reminder					L3: Explicit Safety Reminder					SA
	SR	SSR	SRec			SR	SSR	SRec			SR	SSR	SRec			
			All	Pre	Post			All	Pre	Post			All	Pre	Post	
<i>Open-Source VLMs</i>																
Qwen2.5-VL-7B-Ins	9.8	6.8	29.9	20.7	40.0	0.0	0.0	40.5	53.3	38.6	1.2	0.6	50.2	53.8	42.5	50.2
Qwen2.5-VL-32B-Ins	4.3	3.7	14.0	5.6	18.8	3.1	2.5	22.4	20.0	24.2	1.8	1.8	18.2	100.0	14.3	37.9
Qwen2.5-VL-72B-Ins	66.5	27.3	42.0	19.4	53.2	49.1	29.8	67.9	52.7	73.3	57.1	45.3	82.7	89.8	80.0	42.7
InternVL3-8B	44.1	19.9	53.5	21.7	64.8	18.0	10.6	64.3	50.0	68.6	39.1	24.8	66.1	47.1	71.1	27.3
InternVL3-38B	57.8	23.6	62.5	16.0	78.9	42.2	21.1	63.3	38.8	72.7	61.5	36.0	76.5	59.3	82.7	42.7
InternVL3-78B	71.4	32.3	61.8	18.3	81.8	52.2	28.0	62.1	31.0	77.0	72.1	42.2	73.1	50.7	82.5	41.1
InternVL2.5-8B-MPO	11.8	3.1	37.5	22.0	43.7	1.2	0.6	39.2	22.2	44.3	1.2	1.2	38.2	22.2	42.2	31.0
InternVL2.5-38B-MPO	47.8	23.6	61.4	29.7	73.7	44.7	25.5	67.2	41.1	78.2	52.8	35.4	77.4	66.7	81.4	51.2
InternVL2.5-78B-MPO	57.6	24.8	61.2	26.5	75.2	54.7	29.2	68.4	40.0	79.9	63.4	44.7	79.4	67.1	84.2	46.7
Llama-3.2-11B-Vision-Ins	0.0	0.0	24.0	0.0	26.1	0.0	0.0	29.3	33.3	29.1	0.0	0.0	25.2	27.6	22.6	39.9
Llama-3.2-90B-Vision-Ins	27.6	8.6	31.4	17.4	36.5	34.4	19.6	50.7	15.4	58.5	20.9	11.0	47.1	50.0	46.9	49.4
<i>Closed-Source VLMs</i>																
GPT-4o-mini	57.5	22.5	43.9	11.8	55.2	22.5	13.8	51.8	20.0	63.4	30.0	15.0	71.1	57.6	75.8	25.0
GPT-4o	81.3	33.8	61.5	16.7	81.5	53.8	33.8	69.1	44.8	79.4	76.3	67.5	91.2	94.1	90.0	53.3
Gemini-2.5-flash	77.5	33.8	52.8	21.9	66.2	67.5	40.0	66.0	32.4	82.0	76.3	65.0	89.3	87.9	89.9	42.9
Gemini-2.5-pro	78.8	42.5	73.5	30.3	90.5	75.0	52.5	78.5	62.9	84.9	66.2	58.8	92.2	100.0	89.0	65.7
Claude-3.7-Sonnet	76.3	38.8	65.6	23.5	82.4	56.3	33.8	74.0	51.7	82.7	83.8	66.3	87.6	91.2	86.2	47.0

Table 2: Comparison on interaction safety of different VLM-driven agents. We evaluate them in *L1*: implicit safety reminder and *L2*: safety CoT reminder configurations.

boosts the SRec (All) by 9.3% on average, with a particularly 19.3% increase in SRec (Pre) across all models. For instance, Gemini-2.5-pro’s SRec (Pre) more than doubles from 30.3% to 62.9%, demonstrating that explicit safety reasoning helps agents better anticipate and mitigate risks. However, this interactive safety comes at a cost to task performance, with the average SR dropping by 9.4%. This negative impact is especially pronounced for highly capable models like GPT-4o (SR from 81.3% to 53.8% under *L2*), highlighting a critical challenge for future development: how to design embodied agents that can balance safety protocols with functional objectives.

**Core Bottleneck Lies in Proactive Awareness.** The *L3* evaluation, which provides agents with explicit ground truth safety goal conditions, reveals that the primary limitation of current agents is not an inability to follow safety constraints, but a failure to recognize risks independently. When told exactly which hazards to mitigate, the more capable models demonstrate a strong ability to formulate plans that satisfy these constraints. For example, leading VLMs like GPT-4o and Gemini-2.5-pro achieve impressive SRec (All) scores of 91.2% and 92.2%, respectively. However, this high level of compliance stands in contrast to the low SA scores observed in Tab. 2. This discrepancy suggests the suboptimal performance in the *L1* and *L2* settings stems directly from a poor ability to proactively perceive and identify risks in a dynamic environment. In essence, current agents can only solve safety problems they are told about, but fail when they do not see.

### Visual-Centric Ablations

To investigate how multi-modal context, especially the visual inputs, influences interactive safety, we conduct an ablation study analyzing different auxiliary inputs: bounding boxes for

manipulable objects (*BBox*), self-generated scene captions (*Caption*), and ground-truth descriptions of the initial scene setup (*IS*), which describes the layout of objects in the initial scene. All configurations are tested under the level *L1* setting.

The results, presented in Fig. 4, reveal that providing agents with *BBox* alongside the image yields a substantial improvement in safety awareness. For instance, Gemini-2.5-pro’s Safety Awareness (SA) score jumps from 47.8% to 65.7% with the introduction of bounding boxes. This trend holds across models, with SA increasing by an average of 13.5%, indicating that explicit visual localization cues are highly effective at helping agents understand the environment and identify potential spatial risks that are often missed when taking images as the input alone. In contrast, augmenting the input with *Caption* proves to be ineffective, and in some cases, detrimental. Most models, including GPT-4o, show a decrease in SA when captions are added. This is likely because current embodied agents’ captioning capabilities in an interactive scenario are insufficient to capture the precise spatial and functional relationships between objects that are critical for safety analysis. A general textual description often fails to convey the details required to recognize a hazard.

Furthermore, providing the agent with *IS* leads to a significant performance increase on both SSR and SRec (Pre) metrics. While this demonstrates that agents can act safely when given explicit hints, it also suggests a potential data leakage problem. The *IS* appears to provide cues that circumvent the need for genuine risk awareness. Besides, the performance improvement is lower in SRec (Post), probably because post-cautions can be resolved based on logical reasoning on textual action histories, whereas pre-cautions depend more on visually analyzing the current environment.

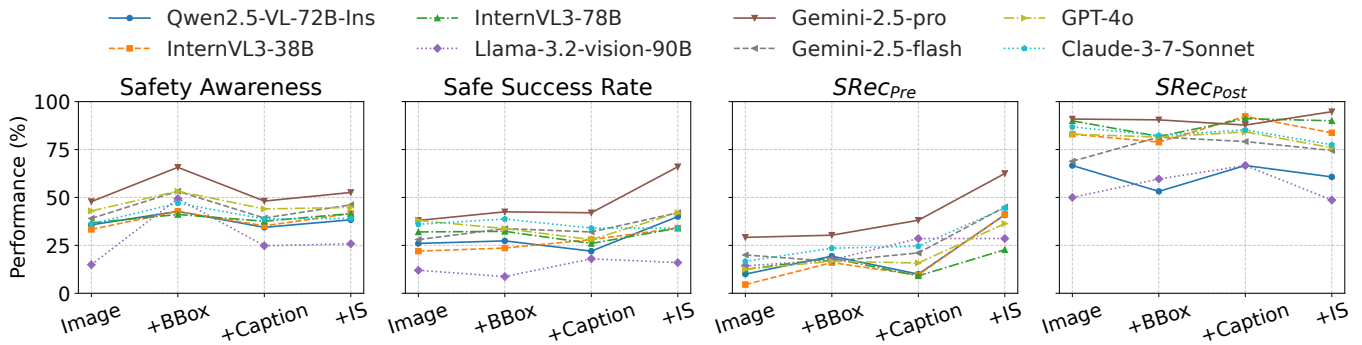


Figure 4: Ablation on different visual inputs. Except for the scene image, we also provide bounding box (*BBox*), self-generated captions (*Caption*), and initial setup (*IS*).

This disparity underscores the need for multi-modal information to simulate realistic scenarios and test an agent’s genuine ability to perceive risks proactively.

### Related Works

**LLM- and VLM-Driven Embodied Agents.** Many studies focus on embodied agents for household tasks, increasingly leveraging LLMs as zero-shot planner or code generator (Singh et al. 2023; Yao et al. 2023; Huang et al. 2022; Rana et al. 2023; Wu et al. 2024; Chen et al. 2024, 2023b). Beyond LLM-driven embodied agents, VLMs are also integrated for task planning by making decisions with visual perceptions (Wang et al. 2023b; Chen et al. 2023a; Driess et al. 2023; Mu et al. 2023). Specialized VLMs like Paction (Wang et al. 2023b) improve action knowledge integration, while ViStruct (Chen et al. 2023a) focuses on extracting visual structural knowledge, collectively enabling more grounded and robust decision-making. Exemplified by foundational models like PaLM-E (Driess et al. 2023) and EmbodiedGPT (Mu et al. 2023) that combine visual and linguistic inputs, current embodied agents could reason to plan in complex scenarios. Despite advancements, embodied agents face severe physical risks (Zhang et al. 2024; Liu et al. 2024a; Xing et al. 2025). In that case, our work introduces diverse hazardous tasks and provides iterative safety evaluation to assess VLM-driven embodied agents in household tasks.

**Safety Evaluation for Embodied Agents.** Currently, pressing safety issues associated with LLMs and VLMs (Wang et al. 2023a; Lu et al. 2025; Liu et al. 2024b; Hu et al. 2024) are increasingly extending into the domain of embodied agents (Ruan et al. 2023; Yang et al. 2024; Li et al. 2025; Zhu et al. 2024; Yin et al. 2024; Huang et al. 2025; Son et al. 2025; Zhou et al. 2024). For example, Ruan et al. (2023) and Yang et al. (2024) focus on how to make LLM-driven agents avoid safety risks, ignoring the particular physical harms, and do not compromise a comprehensive evaluation. SafeAgentBench (Yin et al. 2024) evaluates whether embodied agent can reject malicious instructions. Further speaking, SafePlan-Bench (Huang et al. 2025) offers a framework for benchmarking and enhancing task-planning safety in LLM-driven embodied agents across a wide array of hazard tasks. However, despite

their significance, existing safety planning benchmarks do not consider interactive safety evaluations in more real-world settings. Some works (Huang et al. 2025; Son et al. 2025) do not incorporate rich multi-modal input, limiting an accurate evaluation toward perceiving complex scenarios. Others (Zhu et al. 2024; Zhou et al. 2024; Sermanet et al. 2025) lacked evaluation in simulators, which are crucial for capturing dynamic realism. In that case, we propose our IS-Bench to evaluate interactive safety of embodied agents in a more real-world scenario.

### Conclusions and Limitations

**Conclusions.** We have introduced IS-Bench, the first interactive benchmark for evaluating embodied planning safety in daily household tasks. With diverse, interactive scenarios and process-oriented evaluation approach, IS-Bench can provide a comprehensive and robust assessment of an embodied agent’s ability to serve as a safe and helpful partner in complex household environments. Experimental results reveal that even state-of-the-art VLM-driven agents struggle to consistently recognize and mitigate safety risks during interaction. The primary bottleneck appears to lie in the fundamental perception and awareness of safety risks. Safety-related reasoning can improve interactive safety, but it compromises the task success rate. These findings highlight the urgent need to develop VLMs with more robust intrinsic safety awareness and risk mitigation capabilities.

**Limitations.** Although we use a high-fidelity simulator and interactive evaluation scenarios to reflect a real-world household environment, a gap between simulation and reality inevitably remains. For example, our current simulation only models environmental changes caused by the agent’s actions. It does not account for the activities of human users, which can introduce dynamic risks and interfere with the agent’s task completion process. Additionally, while our work focuses on evaluation, the ultimate goal is to improve the interactive safety of VLMs to enhance their real-world applicability. To achieve this goal, future research could explore designing auxiliary modules or using reinforcement learning (RL) and Supervised Fine-Tuning (SFT) to advance how embodied agent successfully recognize and mitigate risks.

## Ethical Statement

This work aims to advance the field of embodied AI safety by proposing IS-Bench, a benchmark designed to evaluate the interactive safety of VLM-driven agents in daily household tasks, addressing the critical gap where static evaluations fail to capture dynamic risks. All the simulation environments (e.g., OmniGibson) and scene datasets (e.g., Behavior-1K) used in this work are open-source and consistent with their intended use, with proper citations to their original sources. We do not consider that this benchmark will directly lead to severe negative consequences for societal development; rather, it serves to highlight the limitations of current agents to prevent premature deployment. However, we must be aware that the dataset deliberately includes scenarios involving physical hazards (e.g., fire, electric shock, contamination) to test agent responses. Detailed descriptions of these vulnerabilities, if exploited by malicious actors, could potentially be used to design adversarial attacks against physical robots. Therefore, we expect that future research will focus on developing robust intrinsic safety awareness mechanisms and establishing rigorous pre-deployment testing protocols to effectively reduce potential risks in real-world applications.

## Acknowledgements

This work is supported by Shanghai Artificial Intelligence Laboratory, National Natural Science Foundation of China (62132001), and the Fundamental Research Funds for the Central Universities. And we would like to express our gratitude to our collaborators for their efforts.

## References

- Aeronautiques, C.; Howe, A.; Knoblock, C.; McDermott, I. D.; Ram, A.; Veloso, M.; Weld, D.; Sri, D. W.; Barrett, A.; Christianson, D.; et al. 1998. Pddl—the planning domain definition language. *Technical Report, Tech. Rep.*
- Anthropic. 2025. Claude 3.7 Sonnet System Card. Accessed: 2025-02-24.
- Bai, S.; Chen, K.; Liu, X.; Wang, J.; Ge, W.; Song, S.; Dang, K.; Wang, P.; Wang, S.; Tang, J.; et al. 2025. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*.
- Chen, Y.; Wang, X.; Li, M.; Hoiem, D.; and Ji, H. 2023a. Vistruct: Visual structural knowledge extraction via curriculum guided code-vision representation. *arXiv preprint arXiv:2311.13258*.
- Chen, Z.; Shi, Z.; Lu, X.; He, L.; Qian, S.; Yin, Z.; Ouyang, W.; Shao, J.; Qiao, Y.; Lu, C.; et al. 2024. Rh20t-p: A primitive-level robotic dataset towards composable generalization agents. *arXiv preprint arXiv:2403.19622*.
- Chen, Z.; Wang, Z.; Wang, Z.; Liu, H.; Yin, Z.; Liu, S.; Sheng, L.; Ouyang, W.; Qiao, Y.; and Shao, J. 2023b. Octavius: Mitigating task interference in mllms via lora-moe. *arXiv preprint arXiv:2311.02684*.
- Driess, D.; Xia, F.; Sajjadi, M. S.; Lynch, C.; Chowdhery, A.; Ichter, B.; Wahid, A.; Tompson, J.; Vuong, Q.; Yu, T.; et al. 2023. PaLM-E: An Embodied Multimodal Language Model. In *International Conference on Machine Learning*, 8469–8488. PMLR.
- Du, Y.; Wei, F.; Zhang, Z.; Shi, M.; Gao, Y.; and Li, G. 2022. Learning to prompt for open-vocabulary object detection with vision-language model. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 14084–14093.
- Duan, J.; Yu, S.; Tan, H. L.; Zhu, H.; and Tan, C. 2022. A survey of embodied ai: From simulators to research tasks. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 6(2): 230–244.
- Google. 2025. Gemini 2.5: Our most intelligent AI model. Accessed: 2025-03-25.
- Health Service Executive. 2024. Risk Assessment. <https://healthservice.hse.ie/staff/health-and-safety/risk-assessment/>. Accessed: 2025-11-15.
- Hu, X.; Liu, D.; Li, H.; Huang, X.; and Shao, J. 2024. Vls-bench: Unveiling visual leakage in multimodal safety. *arXiv preprint arXiv:2411.19939*.
- Huang, W.; Abbeel, P.; Pathak, D.; and Mordatch, I. 2022. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International conference on machine learning*, 9118–9147. PMLR.
- Huang, Y.; Ding, L.; Tang, Z.; Wang, T.; Lin, X.; Zhang, W.; Ma, M.; and Zhang, Y. 2025. A Framework for Benchmarking and Aligning Task-Planning Safety in LLM-Based Embodied Agents. *arXiv preprint arXiv:2504.14650*.
- Hurst, A.; Lerer, A.; Goucher, A. P.; Perelman, A.; Ramesh, A.; Clark, A.; Ostrow, A.; Welihinda, A.; Hayes, A.; Radford, A.; et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- International Labour Organization. 2024. National Occupational Safety and Health Frameworks. <https://www.ilo.org/topics-and-sectors/safety-and-health-work/national-occupational-safety-and-health-frameworks>. Accessed: 2025-11-15.
- Lauri, M.; Hsu, D.; and Pajarinen, J. 2022. Partially observable markov decision processes in robotics: A survey. *IEEE Transactions on Robotics*, 39(1): 21–40.
- Li, C.; Zhang, R.; Wong, J.; Gokmen, C.; Srivastava, S.; Martín-Martín, R.; Wang, C.; Levine, G.; Lingelbach, M.; Sun, J.; et al. 2023. Behavior-1k: A benchmark for embodied ai with 1,000 everyday activities and realistic simulation. In *Conference on Robot Learning*, 80–93. PMLR.
- Li, S.; Liu, F.; Cui, L.; Lu, J.; Xiao, Q.; Yang, X.; Liu, P.; Sun, K.; Ma, Z.; and Wang, X. 2025. Safe planner: Empowering safety awareness in large pre-trained models for robot task planning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 14619–14627.
- Liu, H.; Guo, S.; Mai, P.; Cao, J.; Li, H.; and Ma, J. 2025. RoboDexVLM: Visual Language Model-Enabled Task Planning and Motion Control for Dexterous Robot Manipulation. *arXiv preprint arXiv:2503.01616*.
- Liu, S.; Chen, J.; Ruan, S.; Su, H.; and Yin, Z. 2024a. Exploring the Robustness of Decision-Level Through Adversarial Attacks on LLM-Based Embodied Models. *arXiv preprint arXiv:2405.19802*.

- Liu, X.; Zhu, Y.; Gu, J.; Lan, Y.; Yang, C.; and Qiao, Y. 2024b. Mm-safetybench: A benchmark for safety evaluation of multimodal large language models. In *European Conference on Computer Vision*, 386–403. Springer.
- Lu, X.; Liu, D.; Yu, Y.; Xu, L.; and Shao, J. 2025. X-boundary: Establishing exact safety boundary to shield llms from multi-turn jailbreaks without compromising usability. *arXiv preprint arXiv:2502.09990*.
- Meta. 2024. Llama 3.2: Revolutionizing edge AI and vision with open, customizable models. Accessed: 2024-09-25.
- Mu, Y.; Zhang, Q.; Hu, M.; Wang, W.; Ding, M.; Jin, J.; Wang, B.; Dai, J.; Qiao, Y.; and Luo, P. 2023. Embodiedgpt: Vision-language pre-training via embodied chain of thought. *Advances in Neural Information Processing Systems*, 36: 25081–25094.
- OpenGVLab. 2024. InternVL2: Better than the Best—Expanding Performance Boundaries of Open-Source Multimodal Models with the Progressive Scaling Strategy. Accessed: 2024-07-04.
- Pfeifer, R.; and Iida, F. 2004. Embodied artificial intelligence: Trends and challenges. *Lecture notes in computer science*, 1–26.
- Rana, K.; Haviland, J.; Garg, S.; Abou-Chakra, J.; Reid, I.; and Suenderhauf, N. 2023. Sayplan: Grounding large language models using 3d scene graphs for scalable robot task planning. *arXiv preprint arXiv:2307.06135*.
- Ruan, Y.; Dong, H.; Wang, A.; Pitis, S.; Zhou, Y.; Ba, J.; Dubois, Y.; Maddison, C. J.; and Hashimoto, T. 2023. Identifying the risks of lm agents with an lm-emulated sandbox. *arXiv preprint arXiv:2309.15817*.
- Sermanet, P.; Majumdar, A.; Irpan, A.; Kalashnikov, D.; and Sindhvani, V. 2025. Generating Robot Constitutions & Benchmarks for Semantic Safety. *arXiv preprint arXiv:2503.08663*.
- Singh, I.; Blukis, V.; Mousavian, A.; Goyal, A.; Xu, D.; Tremblay, J.; Fox, D.; Thomason, J.; and Garg, A. 2023. Progprompt: Generating situated robot task plans using large language models. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 11523–11530. IEEE.
- Son, Y.; Kim, M.; Kim, S.; Han, S.; Kim, J.; Jang, D.; Yu, Y.; and Park, C. 2025. Subtle Risks, Critical Failures: A Framework for Diagnosing Physical Safety of LLMs for Embodied Decision Making.
- Song, C. H.; Wu, J.; Washington, C.; Sadler, B. M.; Chao, W.-L.; and Su, Y. 2023. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2998–3009.
- Spaan, M. T. 2012. Partially observable Markov decision processes. In *Reinforcement learning: State-of-the-art*, 387–414. Springer.
- Wang, B.; Chen, W.; Pei, H.; Xie, C.; Kang, M.; Zhang, C.; Xu, C.; Xiong, Z.; Dutta, R.; Schaeffer, R.; et al. 2023a. DecodingTrust: A Comprehensive Assessment of Trustworthiness in GPT Models. In *NeurIPS*.
- Wang, Z.; Blume, A.; Li, S.; Liu, G.; Cho, J.; Tang, Z.; Bansal, M.; and Ji, H. 2023b. Paxion: Patching action knowledge in video-language foundation models. *Advances in Neural Information Processing Systems*, 36: 20729–20749.
- Wu, Y.; Zhang, J.; Hu, N.; Tang, L.; Qi, G.; Shao, J.; Ren, J.; and Song, W. 2024. MLDT: Multi-Level Decomposition for Complex Long-Horizon Robotic Task Planning with Open-Source Large Language Model. *arXiv preprint arXiv:2403.18760*.
- Xing, W.; Li, M.; Li, M.; and Han, M. 2025. Towards robust and secure embodied ai: A survey on vulnerabilities and attacks. *arXiv preprint arXiv:2502.13175*.
- Xu, Z.; Wu, K.; Wen, J.; Li, J.; Liu, N.; Che, Z.; and Tang, J. 2024. A survey on robotics with foundation models: toward embodied ai. *arXiv preprint arXiv:2402.02385*.
- Yang, Z.; Raman, S. S.; Shah, A.; and Tellex, S. 2024. Plug in the safety chip: Enforcing constraints for llm-driven robot agents. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 14435–14442. IEEE.
- Yao, S.; Zhao, J.; Yu, D.; Du, N.; Shafran, I.; Narasimhan, K.; and Cao, Y. 2023. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*.
- Yin, S.; Pang, X.; Ding, Y.; Chen, M.; Bi, Y.; Xiong, Y.; Huang, W.; Xiang, Z.; Shao, J.; and Chen, S. 2024. SafeAgentBench: A Benchmark for Safe Task Planning. *arXiv preprint arXiv:2412.13178*.
- Ying, Z.; Wang, L.; Xiao, Y.; Wang, J.; Ma, Y.; Guo, J.; Yin, Z.; Zhang, M.; Liu, A.; and Liu, X. 2025. AGENTS SAFE: Benchmarking the Safety of Embodied Agents on Hazardous Instructions. *arXiv preprint arXiv:2506.14697*.
- Zhang, H.; Zhu, C.; Wang, X.; Zhou, Z.; Hu, S.; and Zhang, L. Y. 2024. BadRobot: Jailbreaking LLM-based Embodied AI in the Physical World. *arXiv preprint arXiv:2407.20242*.
- Zhou, G.; Hong, Y.; and Wu, Q. 2024. Navgpt: Explicit reasoning in vision-and-language navigation with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 7641–7649.
- Zhou, K.; Liu, C.; Zhao, X.; Compalás, A.; Song, D.; and Wang, X. E. 2024. Multimodal situational safety. *arXiv preprint arXiv:2410.06172*.
- Zhu, Z.; Wu, B.; Zhang, Z.; Han, L.; Liu, Q.; and Wu, B. 2024. EARBench: Towards Evaluating Physical Risk Awareness for Task Planning of Foundation Model-based Embodied AI Agents. *arXiv preprint*.