

CRAB: CROSS-ENVIRONMENT AGENT BENCHMARK FOR MULTIMODAL LANGUAGE MODEL AGENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

The development of autonomous agents increasingly relies on Multimodal Language Models (MLMs) to perform tasks described in natural language with GUI environments, such as websites, desktop computers, or mobile phones. Existing benchmarks for MLM agents in interactive environments are limited by their focus on a single environment, lack of detailed and generalized evaluation methods, and the complexities of constructing tasks and evaluators. To overcome these limitations, we introduce CRAB, the first agent benchmark framework designed to support cross-environment tasks, incorporating a graph-based fine-grained evaluation method and an efficient mechanism for task and evaluator construction. Our framework supports multiple devices and can be easily extended to any environment with a Python interface. Leveraging CRAB, we developed a cross-platform Crab Benchmark-v0 comprising 120 tasks in computer desktop and mobile phone environments. We evaluated four advanced MLMs using different single and multi-agent system configurations on this benchmark. The experimental results demonstrate that the single agent with GPT-4o achieves the best completion ratio of 38.01%.

1 INTRODUCTION

The development of autonomous agents for human-centric interactive systems—such as desktop OS (Zhang et al., a), websites (Zhou et al.; Koh et al.), smartphones (Zhang et al., b; Xing et al.), and games (Vinyals et al.; Wang et al., a)—has long been an important goal of AI research, aiming to convert natural language instructions into concrete operations. Traditionally, these challenges have been addressed using reinforcement learning (Mnih et al.). Recently, Large Language Models (LLMs) have demonstrated remarkable proficiency in natural language understanding and commonsense reasoning, making them vital tools for developing autonomous agents. This utility is further enhanced by Multimodal Language Models (MLMs), which improve the ability to interpret visual information from GUIs (Cheng et al.).

To effectively develop MLM-based autonomous agents for real-world applications, it is essential to create suitable benchmarks for standardized performance evaluation. However, existing benchmarks still have limitations in terms of interaction methods, platform diversity, evaluation metrics, static task dataset that prevent them from closely mirroring complex real-world applications. First, existing benchmarks that interact with the environments through pre-collected observation data from system environments (Sun et al.; Mialon et al.; Deng et al., 2023) fail to capture the dynamic nature of real-world scenarios without interactive exploration where data and conditions can change unpredictably. Second, existing benchmarks are typically evaluated on a single platform, either Web, Android, or Desktop OS (Shi et al., 2017; Xing et al.; Xie et al.). However, the practical applications usually involve tasks that span multiple platforms. For example, using a smartphone to take a photo and sending it to a desktop for editing with a graphics editor is a common real-world task across multiple platforms. Third, existing evaluation methods are generally either goal-based or trajectory-based (Shi et al., 2017; Xing et al.). Goal-based methods typically employ a coarse-grained binary reward, solely evaluating whether the final system state aligns with the task’s objectives. In contrast, trajectory-based methods can offer more nuanced metrics by assessing the agent’s actions against a gold trajectory yet ignore the possibility of multiple valid pathways to complete a task, making the evaluation results less fair. Lastly, task creation within these complex systems are not static and extensible with fixed templates (Sun et al.; Xie et al.), which limits the diversity and scope of tasks.

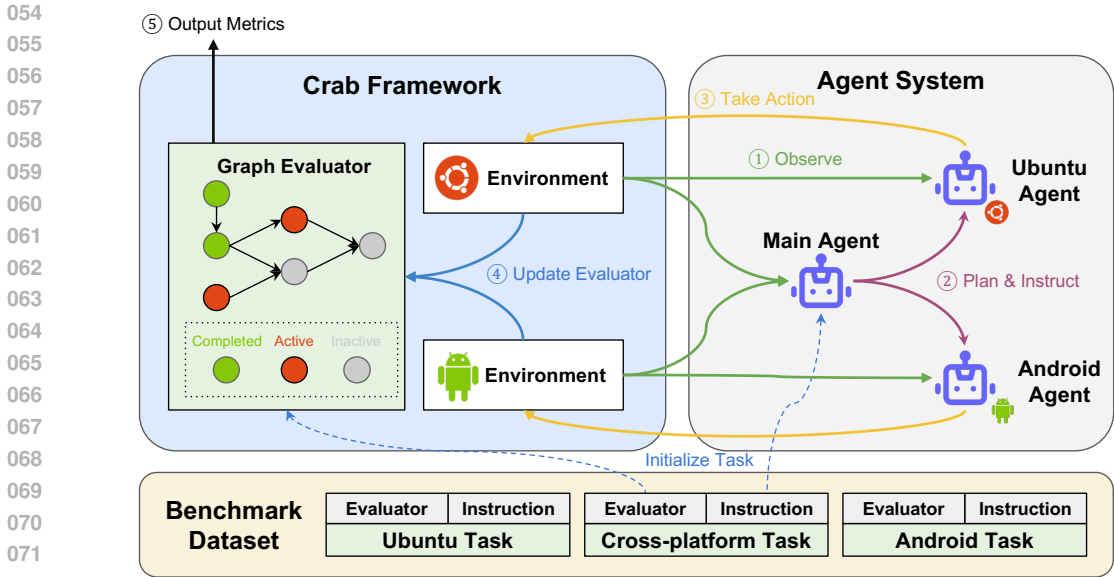


Figure 1: Architecture of the Crab Framework demonstrating a benchmarking workflow for a multi-agent system. A task is initialized by assigning instructions to the main agent and a graph evaluator inside the benchmark system. The workflow progresses through a cycle where the main agent observes, plans, and instructs the sub-agents, who then execute actions within their respective environments. The graph evaluator monitors the status of tasks within the workflow, continuously updating and outputting the task completion metrics throughout the workflow.

We propose a benchmark that closely mirrors real-world situations and an evaluation method that more accurately reflects an agent’s performance on complex tasks. To this end, we introduce CRAB, a novel **C**ross-environment **A**gent **B**enchmark framework. CRAB provides a comprehensive framework for evaluating cross-environment tasks in interactive environments, where the agent needs to operate simultaneously across various devices and platforms, adapting to varied system conditions to complete tasks efficiently. To the best of our knowledge, CRAB is the first autonomous agent benchmark framework that incorporates the **cross-environment tasks**. Moreover, we propose a novel evaluation method called **graph evaluator**. Unlike traditional goal-based and trajectory-based evaluation, our graph evaluator checks the intermediate procedures of completing a task by decomposing the task into multiple sub-goals. Each sub-goal is assigned a judge function to verify its completeness, and each is considered a node in the graph evaluator. The graph structure describes the sequential and parallel relationships between the sub-goals. Therefore, it offers fine-grained metrics similar to trajectory-based evaluations while accommodating multiple valid pathways to a solution, making it more suitable for evaluating tasks that involve various correct approaches. To solve the increasing complexity in cross-environment task construction. We also propose a highly extensible graph-based task construction method called **sub-task composition**. Combining multiple sub-tasks in a graph with task targets allows for efficient construction of various cross-environment tasks with corresponding graph evaluators. The whole framework is implemented in Python and use the network to interact with environments, ensuring easy adaptation to any platform, device, or modality. Table 1 compares CRAB with existing agent benchmark frameworks.

Based on CRAB framework, we propose a benchmark Crab Benchmark-v0 with two cooperated environments that include an Android emulator and an Ubuntu desktop virtual machine. We have developed a total of 120 real-world tasks. These tasks address a wide array of common real-world applications and tools, including but not limited to calendars, email, maps, web browsers, and terminals, and facilitate common interactions between smartphones and desktops. Considerable time has been invested in verifying the accuracy and comprehensiveness of the instructions for subtasks, as well as the generalization and correctness of their evaluators. Most tasks are constructed using a careful composition of sub-tasks, while some tasks are crafted manually to accommodate specific collaborative scenarios. We test 6 popular MLMs, including GPT-4 Turbo, GPT-4o, Claude 3 Pro, Gemini 1.5 Pro, Pixtral-8B, and LLaVA-OneVision-72B across different structures of single-agent

Table 1: **Comparison of existing agent benchmark frameworks.** The columns details key features of each framework: *Interactive Environment* indicates the presence of either interactive environments or static datasets; *Multimodal Observation* specifies the availability of vision-based observations (e.g. screenshots); *Cross-platform* denotes support for multiple operating systems or platforms; *Evaluation* describes the evaluation metrics, categorized as *Goal-based* (checking environment state according solely on the final goal), *Trajectory-based* (comparing agent action trajectory with a gold actions sequence), *Multiple* (varied across tasks), or *Graph-based* (a DAG with each node as an intermediate checkpoint); *Task Construction* shows the task construction method, including *Handmade* (handcrafted by human), *LLM-inspired* (using LLM to generate task drafts but still verified and annotated by human), *Template* (generated by filling in the blanks in task templates), or *Sub-task Composition* (composing multiple sub-tasks to construct tasks and evaluators).

	Interactive Environment	Multimodal Observation	Cross-platform	Evaluation	Task Construction	# of apps or websites
MINIWOB++ (Shi et al., 2017)	Web	✓	✗	Goal-based	Handmade	1
WEBSHOP (Yao et al., 2022)	Web	✓	✗	Goal-based	Template	1
METAGUI (Sun et al.)	✗	✗	✗	Trajectory-based	Handmade	6
GAIA (Mialon et al.)	✗	✗	✗	Goal-based	Handmade	n/a
MIND2WEB (Deng et al., 2023)	✗	✗	✗	Goal-based	LLM-inspired	137
AGENTBENCH (Liu et al., 2024)	Multi-isolated	✗	✗	Multiple	Handmade	n/a
INTERCODE (Yang et al., b)	Code	✗	✗	Goal-based	Handmade	n/a
WEBARENA (Zhou et al.)	Web	✓	✗	Goal-based	Template	6
OMNIACT (Kapoor et al.)	✗	✗	✗	Trajectory-based	Handmade	60+
VWEBARENA (Koh et al.)	Web	✓	✗	Goal-based	Template	4
ANDROIDARENA (Xing et al.)	Android	✓	✗	Trajectory-based	LLM-inspired	9
OSWORLD (Xie et al.)	Desktop OS	✓	✗	Goal-based	Template	9
OSWORLD (Xie et al.)	Desktop OS	✓	✗	Goal-based	Template	9
ANDROIDWORLD (Rawles et al., 2024)	Android	✓	✗	Trajectory-based	Handmade	20
WAA (Bonatti et al., 2024)	Desktop OS	✓	✗	Trajectory-based	Handmade	6
CRAB	Desktop OS & Android	✓	✓	Graph-based	Sub-task Composition	23

and multi-agent systems, totaling 12 different agent settings in our benchmarks. The experimental results show that the single agent structure with GPT-4o model achieves the best overall completion ratio of 38.01%, underscoring the necessity for ongoing development of more effective autonomous agents. Our proposed metrics successfully distinguish between different methods better than previous metrics. We further analyze the different termination reasons that reflect the problems inherent in the communication within the multi-agent system.

2 RELATED WORK

Leveraging LLMs as reasoning units has become an effective approach (Wang et al., 2024b; Huang et al., 2022; Xi et al.) for building autonomous agents, including embodied agents (Wang et al., a; Song et al., 2023; Chen et al., 2023), social simulations (Park et al., 2023; Lin et al., 2023), web navigation (Lù et al.), game playing (Lan et al., 2023; Tan et al., 2024), office assistants Li et al. (2024b), and code generation (Zhang et al., 2023). Specifically, some works apply LLMs to the planning of embodied agents in complex environments (Wang et al., a; Song et al., 2023; Chen et al., 2023). Others focus on simulating human behaviors and social communication by harnessing LLMs’ remarkable human-like understanding and generation capabilities (Park et al., 2023; Lin et al., 2023). Additionally, multi-agent systems have been introduced to enhance the simulation of human behavior through agent cooperation (Li et al., 2023; Hong et al., 2023; Wu et al., 2023; Jin et al., 2024; Wang et al., 2024a). In another approach, several studies have expanded the capacities of agents by incorporating multimodal understanding, enabling agents to process diverse modalities of input data such as images and text (Hong et al.; Liu et al., a; Furuta et al., 2024; Chen et al., 2024).

Various benchmarks are developed to validate the performance of autonomous agents based on the reproducible environments. Miniwob++ (Shi et al., 2017) analyzes the open-domain web tasks, builds corresponding web environment, and produces high-quality datasets considering extensive website and operation categories. GAIA (Mialon et al.) proposes a benchmark which considers the challenges of emergency cases. Mind2Web (Deng et al., 2023) proposes a benchmark for the real-world websites which are genuine and unpredictable, with a high coverage of domains, websites, tasks, and user-interactions. WebArena (Zhou et al.) provides a realistic and reproducible web environment to simulate sufficiently complex web tasks. Several works (Koh et al.; He et al., 2024) further broaden the web tasks, considering the visual tasks to build the benchmark for multi-

modal autonomous agents. SWEBench (Jimenez et al.) builds a benchmark based on the Github, focusing on the coding capacity of understanding and solving github issues. AgentBench (Liu et al., 2024), significantly expands the scope of agent applications within the domain of computer interaction tasks. This expansion is particularly noteworthy as it encompasses the examination of these tasks across a diverse array of complex and challenging environments. OMNIACT (Kapoor et al.) incorporates the visual information of OS screen UI via segmentation and corresponding tagging, which creates corresponding tasks upon the basic elements. OSWorld(Xie et al.) pays attention to the simulations across diverse computer systems, taking XML and screenshots as both inputs and meticulously delineating a standardized format for both the environment and the evaluation process. WindowsAgentArena(Bonatti et al., 2024) focuses on the simulation of windows environment, proposes a challenging set of windows-oriented task, gives a trustful evaluation for the popular environment.

Contemporary studies not only focus on tasks related to control within web and computer systems but also extend their scope to encompass control tasks within mobile systems. MetaGUI (Sun et al.) divides the mobile system control tasks into dialogues and GUI operation traces, collecting GUI traces based on the collected dialogues. AITW (Rawles et al., 2023) produces a large dataset upon a large dataset of real-world scenarios, and builds challenging multi-steps tasks based on the annotated single-step tasks as a two-stage manner. MobileAgent (Wang et al., b) proposes tasks based on Ant Intelligent Assistant(AIA) system, which integrates Standing Operating Procedure(SOP) information for the creation of subtasks. AITZ (Zhang et al., 2024) constructs datasets with Chain-of-Thought (COT) considerations, adding semantic annotations according to visual models at each step, and developing the operational procedure for selected tasks. Mobile Agent BenchWang et al. (2024c) collects app event signals via android accessibility service, builds the benchmark with well annotated operation trajectories, and divides the tasks into several levels. Android World(Rawles et al., 2024) establishes a fully functional environment for the Android system and provides a robust and reliable evaluation of the agent’s capacity in Android-oriented tasks. Although these benchmarks try to evaluate the capacity of agent in a wide range of applications, these benchmarks are built on human annotated trajectories, which are lack of scalability. In the other hand, the evaluation methods of previous benchmarks are highly related to annotated trajectories or goals, which makes the evaluation biased and unfaired.

3 DEFINITIONS

3.1 PROBLEM FORMULATION

Consider autonomous agents performing a task on a digital device (i.e. desktop computer). Such a device typically has input devices (i.e. mouse and keyboard) for human interaction and output devices (i.e. screen) to allow human observation of its state. In CRAB, we represent this type of device as an **environment**. Formally, this environment is defined as a reward-free Partially Observable Markov Decision Process (POMDP), denoted by the tuple $M := (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{O})$, where \mathcal{S} represents the state space, \mathcal{A} the action space, $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ the transition function, and \mathcal{O} the observation space. Considering the collaborative nature of multiple devices in real-world scenarios, we can combine multiple environments into a set $\mathbf{M} = M_1, M_2, \dots, M_n$, where n is the number of environments and each environment $M_j = (\mathcal{S}_j, \mathcal{A}_j, \mathcal{T}_j, \mathcal{O}_j)$. We define a task that requires operations across multiple environments as a **cross-environment task**. This task is formalized as a tuple (\mathbf{M}, I, R) , in which \mathbf{M} is the environment set, I is the task objective in the form of natural language instructions, and R is the reward function of the task. An **agent system**, designed to complete a task represented by an instruction I , can be modeled as a policy $\pi((m, a) \mid (I, H, o_1, \dots, o_n))$, which defines the probability of taking action a in environment m when receiving observation (o_1, \dots, o_n) from environment (M_1, \dots, M_n) with a history of actions H . An **agent** within the agent system operates with a fixed back-end MLM, a predefined system prompt, and retains its chat history. An agent system is composed of either a single agent responsible for all planning, reasoning, and action-taking or multiple agents connected through a communication workflow to collaborate.

3.2 GRAPH OF TASK DECOMPOSITION

Decomposing a complex task into several simpler sub-tasks has been proved to be an effective prompting method for LLMs (Khot et al., 2023). Some studies repre-

216 sent sub-tasks in a graph structure. For instance, PLaG (Lin et al.) uses a graph-
 217 based structure to enhance plan reasoning within LLMs, while DyVal (Zhu et al.,
 218 2024) employs directed acyclic graphs (DAGs) to facilitate dynamic evaluation of LLMs.
 219 By introducing this concept into the realm
 220 of benchmarks, naturally, decomposing
 221 a complex task into sub-tasks that have
 222 both sequential and parallel connections
 223 forms a DAG. Therefore, we introduce
 224 the **Graph of Decomposed Tasks (GDT)**,
 225 which provides a new task decomposition
 226 method representing decomposed sub-
 227 tasks within a DAG structure. In GDT, each
 228 node is a sub-task, formalized as a tuple
 229 (m, i, r) , where m specifies the environ-
 230 ment in which the sub-task is performed, i
 231 provides the natural language instruction,
 232 and r represents the reward function. This
 233 function evaluates the state of m and out-
 234 puts a boolean value to determine if the
 235 sub-task is completed. The edges within
 236 GDT represent the sequential relationship
 237 between sub-tasks. An example GDT is
 238 shown in Fig. 2.

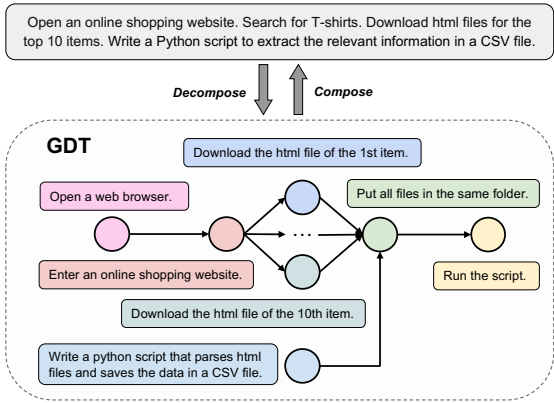


Figure 2: An example of a Graph of Task Decomposition.

239 4 THE CRAB FRAMEWORK

241 4.1 CROSS-ENVIRONMENT AGENT INTERACTION

243 Compared to single-environment tasks, cross-environment tasks offer two main advantages for
 244 benchmarking agents. First, cross-environment tasks reflect real-world scenarios where humans
 245 use multiple devices simultaneously to accomplish tasks. Second, these tasks require sophisticated
 246 message processing and information transfer between environments. Such tasks demand that the agent
 247 plan actions, construct outputs for each environment, and remember what needs to be transferred,
 248 showcasing a high-level understanding of real-world and ability to solving complex tasks. CRAB
 249 uses a unified interface for agents to operate in all environments. We define an action by its name,
 250 the environment it belongs to, a concrete description of its functionality, and the parameters with
 251 descriptions. The agent must provide the action name, parameters, and the target environment in
 252 each turn. CRAB translates the action into its corresponding function and routes it to the physical or
 253 virtual device through the network. Through this approach, CRAB can easily adapt to any platform
 254 in any modality, whether it’s a physical or virtual device, or even a single application like a browser,
 255 by simply defining several interactive functions. Implementation details are in the Appendix A.2.

256 4.2 GRAPH EVALUATOR

258 To assess the capabilities of MLM agents, most benchmarks (Shi et al., 2017; Deng et al., 2023; Koh
 259 et al.; Zhou et al.) evaluate agents based on solely the final states of the environment after agent
 260 operations. Typically, they only judge whether the final goal is success or fail. However, this approach
 261 does not capture incremental progress made by the agents. For instance, consider two agents tasked
 262 with installing a new application on a computer: agent a successfully downloads the installer but fails
 263 during the installation process, whereas agent b does not even try to find the installer. Despite Agent
 264 a making more progress, both are deemed failures under the goal-based evaluation system, resulting
 265 in an unfair assessment of their performance. An alternative method, *Trajectory-based Matching*
 266 (Xing et al.; Kapoor et al.), abandons state-based evaluation and instead compares the agent’s actions
 267 against a predefined gold action sequence for each task, giving nuanced metrics. Nevertheless, this
 268 method faces challenges in real-world systems where tasks may have multiple valid execution paths.
 269 Inspired by the "decomposing" idea from GDT (Sec. 3.2), we propose a novel integrated approach,
 the *Graph Evaluator*, which provides fine-grained metrics and supports multiple valid paths.

270 To build a graph evaluator for a given task, we begin by decomposing the task into a GDT, where
 271 each sub-task is associated with an intermediate environment state critical to completing the overall
 272 task. Nodes in the graph evaluator activate when they either have no incoming edges or after all
 273 their preceding tasks are completed, ensuring a sequential order of tasks. After an agent takes an
 274 action, the system checks these active nodes to verify if the target state of each node is reached. A
 275 node completion triggers successor nodes to activate and verify the state. This cycle repeats until
 276 no new nodes activate, showing that the system’s task sequence aligns with the current state of the
 277 environment. Unlike trajectory-based methods, which compare sequences of agent actions, the Graph
 278 Evaluator does not rely on the specific actions taken by the agent, allowing it the freedom to choose
 279 any path. Instead, it concentrates on the key intermediate states of the environment necessary for
 280 reaching the final goal.

281 4.3 METRICS

283 Given a Graph Evaluator synchronized with the environment state, it becomes possible to track
 284 agent progress through the current status of sub-task completions. Beyond the traditional **Success**
 285 **Rate (SR)**, which marks a task as *success* only when all sub-tasks are completed, we introduce
 286 three metrics aiming at assessing both performance and efficiency of agents, leveraging the detailed
 287 sub-task status provided by the graph evaluator. Specifically, the **Completion Ratio (CR)** measures
 288 the proportion of completed sub-task nodes relative to the total nodes in the graph, calculated as
 289 C / N , where C is the number of completed nodes and N is the total number of nodes. This
 290 metric offers a straightforward measure of an agent’s progress on a given task. The **Execution**
 291 **Efficiency (EE)**, calculated as CR / A , where A denotes the count of executed actions. It evaluates
 292 how efficiently actions are executed relative to the completion of nodes, reflecting the agent’s task
 293 execution efficiency. Lastly, the **Cost Efficiency (CE)**, calculated as CR / T , where T is the total
 294 number of model tokens used, evaluates the efficiency of resource consuming by the agent.

295 4.4 TASK AND EVALUATOR CONSTRUCTION

296 Despite the graph evaluator offers detailed evaluations, one challenge is the complexity in creating
 297 each evaluator. Creating a graph evaluator requires: (1) adequately decomposing a task into multiple
 298 sub-tasks, each with a well-defined graph structure; and (2) engaging an expert who is well-acquainted
 299 with the target platform to carefully craft an evaluator for each sub-task. To streamline the creation of
 300 tasks and the development of evaluators, we consider to build GDTs by sub-tasks.

302 There are two primary challenges in constructing GDT: (1) Sub-tasks still require manual creation,
 303 necessitating a method to quickly generate them on a large scale; (2) Properly modeling the sequential
 304 and parallel relationships between sub-tasks, ensuring that the edges connecting sub-task nodes are
 305 semantically meaningful and systematically applicable. A template-based approach is commonly
 306 used to address the first issue by generating a large number of tasks efficiently. To tackle the second
 307 challenge, we employ the message transferring concept (Sec. 4.1). Specifically, if a sub-task α
 308 produces an output message that serves as an input for another sub-task β , then α can be considered
 309 a legitimate prerequisite of β , allowing us to connect α and β with an directed edge in the GDT. To
 310 further refine our approach, we introduce a *sub-task template* structure. Each sub-task is described
 311 using a natural language instruction template that includes several replaceable input attributes. The
 312 types of each input attribute and the task output should be defined carefully. To generate a GDT,
 313 input attributes can be filled with either a hand-crafted value corresponding to their type or linked to
 314 a task with the same output type as the input type. From the evaluator’s perspective, each sub-task
 315 template is linked to an evaluator generator that uses the input attribute value to generate evaluator
 316 subgraphs. Once a GDT is constructed, the graph evaluator is created by interlinking each subgraph.
 317 The description for the composed task is initially generated by GPT-4 using the sub-task descriptions
 318 as prompts and subsequently refined and polished by human reviewers. Compared to the naive
 319 template-based approach, our sub-task templates are not only more detailed but also designed with
 320 scalability in mind, ensuring that tasks can be composed correctly across various scenarios. Instead of
 321 simple string replacements, our templates are built from sub-task blocks that can be freely combined
 322 with other blocks, allowing for the creation of a wide variety of tasks with descriptions generated by
 323 LLM. Our task generation method ensures that the graph evaluators can be automatically generated,
 relieving users from the need to write any code. This level of automation and ease of use is intended
 to make the system accessible to a broader range of users.

5 THE CRAB BENCHMARK

Environments. We build an agent benchmark Crab Benchmark-v0 featuring with cross-environment, graph evaluator, and task generation through CRAB framework. The environments consists of an Android smartphone emulator and a Ubuntu Linux desktop virtual machine. We establish both environments in a reproducible and standalone manner and utilize snapshots to ensure a consistent initial state for all environments. The observation space consists solely of the current system screen for both environments, captured in image format at each step of the agent’s interaction. We employ the Set-of-Marks visual prompt method (Yang et al., a) to label each interactive element on the screen. Interactive elements are identified using the GroundingDINO (Liu et al., b) with `icon.logo.text` prompt to locate all interactive icons. Additionally, Optical Character Recognition (OCR) is utilized through EasyOCR¹ to detect and label interactive text elements. Each detected item is assigned a unique integer ID, facilitating reference within the action space. The action spaces for Ubuntu and Android are distinct and designed to be close to the common interactions in the real devices. For Ubuntu, we define the following actions: mouse-based actions, keyboard-based actions and a shortcut action to search for applications. For Android, the action set includes tapping actions, a text action, a physical button action, and an action to open the app drawer. Additionally, we introduce three environment-irrelevant actions: completing the task, submitting an answer and waiting. Detailed descriptions for the environment implementation are shown in Appendix A.1.

Tasks. We meticulously construct 17 sub-task templates for the Android environment and 19 sub-task templates for the Ubuntu environment. The Ubuntu templates encompass a variety of tasks such as Command Line Interface (CLI) operations, file system management, search engine usage, desktop configurations, and map navigation. Conversely, the Android sub-task templates are primarily focused on the storage and transmission of messages via various applications. Each sub-task template is linked to a graph evaluator consisting of one to four nodes. Each sub-task are its graph evaluator is verified by at least two related field experts. We make sure that all tasks are reachable by human. We generate 104 tasks by sub-task composition and make 16 tasks by hand to include more complex scenarios that cannot easily described by the sub-tasks. The dataset has 29 Android tasks, 73 Ubuntu tasks and 18 cross-platform tasks, totaling 120 tasks. Our tasks are intentionally designed to be more complex than those in other benchmarks, which naturally requires more time for design and experimentation. A single sub-task in our benchmark might involves multiple operations across several applications, unlike prior works where most tasks often focus on solving problems within a single application. With multiple applications nature combined with the scalability of our task composition and graph evaluator, our tasks are sufficiently challenging to test an agent’s performance across different applications and scenarios, thereby effectively assessing its generalization ability. The format and the applications covered by the dataset are shown in Appendix A.3 and A.4, respectively.

Evaluators. To assess the intermediate states of sub-tasks as described in Sec. 4.2, we have implemented a comprehensive suite of execution-based evaluators. These evaluators retrieve and assess specific current states, such as the edited content of a file or a modified setting, thereby determining the successful completion of a sub-task. For each evaluator, input attributes are carefully selected to interpret software information or system settings relevant to the scenario defined for the sub-task. For instance, evaluators use file paths before and after edits as input parameters to verify the completion of file editing sub-tasks. Specifically, for sub-tasks on the Android platform, we incorporate XML-based evaluators (Xing et al.). We dump UI layout as XML path and verify whether the UI content matches the expected state. For the Ubuntu platform, we employ image matching techniques (Potje et al., 2024; Jiang et al., 2024; Edstedt et al., 2024) and OCR to handle scenarios where acquiring necessary state information through conventional APIs is challenging. Image matching offers fine-grained visual correspondences by comparing keypoint features between images, allowing us to assess spatial relationships among visual elements. Using OCR and image matching, we can accurately evaluate tasks such as verifying whether an agent has successfully created a slide with specified images, text content, and layouts—tasks for which trivial evaluation methods are lacking. We utilize EasyOCR¹ and XFeat² as our primary tools for OCR and image matching. For tasks with real-time characteristics that may change over time, we implement crawler scripts to capture dynamic values at the moment of evaluation. These values are then compared with the results achieved by the

¹<https://github.com/JaidedAI/EasyOCR>

²https://github.com/verlab/accelerated_features

agent upon task completion. We have a total of 59 evaluator functions with different types. Each task has 4.2 evaluators in average of the whole dataset.

6 EXPERIMENTS

6.1 BASELINE AGENT SYSTEM

At the core of MLM Agents are backend Multimodal Language Models that provide natural language and image understanding, basic device knowledge, task planning, and logical reasoning abilities. To run in Crab Benchmark-v0, the backend model needs to support: (1) Accept multimodal mixed input, as the system provides both screenshots and text instructions as prompts; (2) Handle multi-turn conversations, as most tasks require the agent to take multiple actions, necessitating the storage of history messages in its context; (3) Generate structured output through function calling, ensuring the proper use of provided actions with type-correct parameters. However, most open source models do not provide explicit function calling feature, we let these models generate structured JSON output to simulate the function calling behavior.

We selected 4 commercial and 2 open source MLMs that meet these criteria for our experiments: GPT-4o (gpt-4o-2024-05-13) (OpenAI, 2024), GPT-4 Turbo (gpt-4-turbo-2024-04-09) (Achiam et al.), Gemini 1.5 Pro (May 2024 version) (Reid et al.), Claude 3 Opus (claude-3-opus-20240229) (Anthropic, Year), Pixtral-12B (Pixtral-12B-2409)³, and LLaVA-OneVision-72B (llava-onevision-qwen2-72b-ov-chat) (Li et al., 2024a). These models serve as the backend models for our agents. Specifically, We use function calling feature in the four commercial models and JSON output in the two open source models that do not support function calling. Since the JSON output setting uses different prompts from the other, we employ a GPT-4o agent without function calling as the control group to the open source models.

Beyond the MLM backend, the structure of agent systems also influences overall performance. To examine how different multi-agent structures impact performance, we design three agent system structures: **single agent**, **multi-agent by functionality**, and **multi-agent by environment**. In the **single agent** structure, one agent manages all responsibilities, including observation analysis, planning, reasoning, and format the output action. The **multi-agent by functionality** structure splits tasks between a main agent, responsible for analysis and planning, and a tool agent that translates instructions into actions without accessing environmental observations. This division allows the main agent to concentrate on high-level tasks without managing functional call formats. Meanwhile, in the **multi-agent by environment** setup, responsibilities are further distributed. A main agent processes all environmental observations for high-level planning, while each environment-specific sub-agent executes actions based on the main agent’s instructions, incorporating observations from their respective environments.

For all models, we utilized the default API parameters and retained two turns of historical messages to ensure messages do not exceed the context window. The interaction turns are limited to 15 and the task will be terminated when reaching max turns. The agent can also terminate the task ahead if it thinks the task is completed. The screenshots are passed through PNG format with the highest quality that the APIs provide. Detailed agent and prompt designs are shown in Appendix B. In the experiment, we deployed four cloud machines cloned from the same disk image to ensure a consistent environment for all agents. Running a single agent setting in the benchmark requires at least 30 hours to complete on one machine. This duration depends on the API call times and the necessity for manual resets in certain tasks.

6.2 RESULT

The primary outcomes are detailed in Table 2. Aside from the *Success Rate*, *Completion Rate*, *Execution Efficiency*, and *Cost Efficiency* mentioned above, we also present the reasons for agent termination to further investigate the factors preventing the agent system from completing the task.

Comparison of backend models. The GPT-4o and GPT-4 Turbo models, developed by OpenAI, achieved the highest average success rates and completion ratios (CR) among the tested models.

³<https://mistral.ai/news/pixtral-12b/>

Table 2: **Evaluation results on Crab Benchmark-v0.** The *Model* column identifies the backend masked language models (MLMs) used. The *Structure* column describes the configuration of the agent system: *Single* means *single agent*; *By Func* is *multi-agent by functionality*; *By Env* indicates *multi-agent by environment*. We provide traditional metric of *Success Rate* (SR) alongside newly introduced metrics: *Completion Ratio* (CR), *Execution Efficiency* (EE), and *Cost Efficiency* (CE). Note that Gemini 1.5 Pro has an invalid CE because the Gemini API does not support retrieving token counts at the start time of experiments. The *Termination Reason* shows the ratio of reasons why the agent is terminated when the task is not success. *False Completion* (FC) indicates that the agent believes it has completed the task, but it actually has not; *Reach Step Limit* (RSL) means the agent has reached the step limit but has not completed the task; *Invalid Action* (IA) refers to the agent producing outputs that do not follow instructions, which may include invalid formats, nonexistent actions, or invalid action parameters.

Agent system		Metrics				Termination Reason		
Model	Structure	SR(%) ↑	CR(%) ↑	EE(%) ↑	CE(%) ↑	FC(%)	RSL(%)	IA(%)
GPT-4o	Single	14.17	38.01	4.15	5.29×10^{-4}	8.33	55.83	21.67
GPT-4o	By Func	15.00	34.00	3.93	5.31×10^{-4}	10.83	54.17	20.00
GPT-4o	By Env	14.17	33.34	3.84	2.74×10^{-4}	8.33	48.33	29.17
GPT-4 TURBO	Single	9.17	33.35	3.80	4.52×10^{-4}	8.33	65.00	17.50
GPT-4 TURBO	By Func	13.33	33.48	4.07	4.38×10^{-4}	10.83	40.00	35.83
GEMINI 1.5 PRO	Single	5.00	15.48	1.72	n/a	2.50	55.83	36.67
GEMINI 1.5 PRO	By Func	5.00	12.76	1.42	n/a	8.33	33.33	53.33
CLAUDE 3 OPUS	Single	3.33	19.60	1.95	1.85×10^{-4}	10.00	57.50	29.17
CLAUDE 3 OPUS	By Func	3.33	16.48	1.72	1.77×10^{-4}	28.33	34.17	34.17
GPT-4o w/o FC	Single	9.17	23.05	2.34	3.93×10^{-4}	5.00	42.50	43.33
PIXTAL-12B	Single	0.83	9.50	0.75	0.87×10^{-4}	0.83	75.83	22.50
LLAVA-OV-72B	Single	0.83	6.64	0.52	1.02×10^{-4}	12.50	71.67	15.00

Specifically, GPT-4o, with a CR of 38.01%, slightly outperforms GPT-4 Turbo, which has a CR of 33.35%. This result suggests that GPT-4o may have been trained on more GUI-related data. Claude 3 outperforms Gemini 1.5 in terms of CR, but there remains a significant gap between the GPT-4 series and other models. Claude and Gemini have a higher Invalid Action Ratio, usually failing by clicking nonexistent elements on the screen or taking nonexistent actions. Regarding efficiency, the GPT-4 series also demonstrates strong performance, with GPT-4o having a higher CE value compared to GPT-4 Turbo, highlighting its cost-effectiveness. Considering the models of JSON output, GPT-4o’s performance drops compared to the function-calling-enabled version, primarily due to its higher Invalid Action rate. This underscores the effectiveness of function calling in generating structured output. However, GPT-4o still significantly outperforms open-source models. In open source models, Pixtral-12B, with far fewer parameters, achieves a better CR compared to LLaVA-ov-72B, showcasing its efficiency. Although the open-source models generally understand screenshots and generate step-by-step plans correctly, they often fail to execute the correct actions according to the plan. Moreover, they do not effectively analyze task completion through observation. Once an incorrect action is performed, they tend to assume current step is success and proceed to the next step.

Comparison of agent structures. The performance of multi-agent structures on all backend MLMs is slightly lower than that of single-agent structures, which is somewhat unconventional. Based on the communication log, we find that multi-agent structures tend to experience information loss during inter-agent communication, leading to misunderstandings among downstream agents. This increases the likelihood of multi-agent structures taking invalid actions and incorrectly completing tasks. These experiments demonstrate that the design of the communication protocol and selecting the appropriate scenario are crucial for multi-agent systems. A detailed analysis is included in Appendix C.2. In terms of efficiency, multi-agent structures require more chat rounds, which can consume more tokens, resulting in a lower CE compared to single-agent settings.

Comparison of platforms. We have three types of tasks: Ubuntu, Android, and cross-environment. The metrics for each type of task can reveal the model or structure preferences. As shown in Fig. 3,

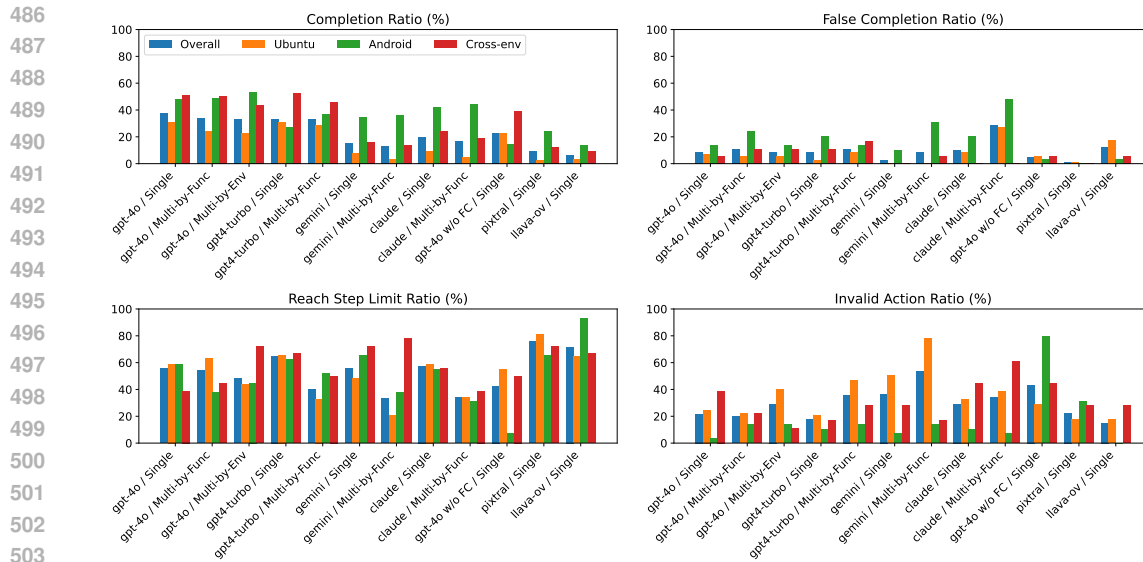


Figure 3: Completion Ratio and Termination Reasons on different platforms.

the GPT-4o model demonstrates significantly better performance on Android and cross-platform tasks compared to GPT-4 Turbo, which may indicate OpenAI’s increased focus on mobile devices. Additionally, models like Gemini, Claude, Pixtral, and LLaVa-OV perform better on Android devices compared to the Ubuntu, likely due to less training on Linux desktop data, which makes it difficult for them to recognize desktop icons. The invalid action ratio is an interesting metric to observe how strictly an agent system follows instructions and outputs valid formats. While it does not fully represent agent performance, it’s notable that the two open-source models exhibit a low invalid action ratio but still fail to complete tasks. By comparing GPT-4o with function calling and using JSON output, the importance of function calling for structured output is emphasized, particularly in the Android environment, where the invalid action ratio is extremely high without it. This phenomenon suggests a hypothesis that the more creative the model is, the less ability it has to follow instructions, resulting in more hallucinations. We include further platform specific results in Appendix C.1.

Comparison of metrics. The completion ratio metric reveals a notable performance difference between models. For instance, even though GPT-4o with single agent structure and with multi-agent by environment structure have the same success rates, their completion ratios differ by up to 4.67%. This highlights the value of the completion ratio in assessing the effectiveness of different methods. For a more detailed analysis of each model and structure’s performance, we provide several case studies in the Appendix. C.3.

7 CONCLUSION

We propose the CRAB framework, which introduces the cross-environment automatic task-performing problem, featuring advanced graph-based task generation and evaluation methods that reduce manual effort in task design and provide more dynamic and accurate agent assessments. Based on this framework, we present Crab Benchmark-v0, a set of high-quality cross-environment tasks in smartphone and desktop environments, equipped with advanced visual prompting techniques. We tested various backend models and agent system structures on this dataset. The results reveal preferences for different agent settings, demonstrating Crab Benchmark-v0’s strong ability to distinguish MLMs and autonomous agent systems. Despite our contribution to advancing cross-environment agent research, there are still some limitations. The sub-tasks are built upon the original apps in the Ubuntu and Android systems on Pixel devices, which limits the coverage of a wider range of applications. The current visual prompting methods do not fully recognize all interactive elements, hindering agent performance. Future work can focus on expanding the dataset and environments, testing more models, prompts, and multi-agent structures, as well as improving the use of visual prompting methods within the benchmark.

REFERENCES

- 540
541
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
543 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
544 URL <http://arxiv.org/abs/2303.08774>.
- 545 Anthropic. The claude 3 model family: Opus, sonnet, haiku. https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf, Year.
- 546
547
548 Fabrice Bellard. Qemu, a fast and portable dynamic translator. In *USENIX annual technical*
549 *conference, FREENIX Track*, volume 41, pp. 10–5555. California, USA, 2005.
- 550
551 Rogerio Bonatti, Dan Zhao, Francesco Bonacci, Dillon Dupont, Sara Abdali, Yinheng Li, Yadong Lu,
552 Justin Wagle, Kazuhito Koishida, Arthur Bucker, Lawrence Jang, and Zack Hui. Windows agent
553 arena: Evaluating multi-modal os agents at scale, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2409.08264)
554 [2409.08264](https://arxiv.org/abs/2409.08264).
- 555 Wentong Chen, Junbo Cui, Jinyi Hu, Yujia Qin, Junjie Fang, Yue Zhao, Chongyi Wang, Jun Liu,
556 Guirong Chen, Yupeng Huo, Yuan Yao, Yankai Lin, Zhiyuan Liu, and Maosong Sun. Guicourse:
557 From general vision language models to versatile gui agents, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2406.11317)
558 [2406.11317](https://arxiv.org/abs/2406.11317).
- 559 Yongchao Chen, Jacob Arkin, Yang Zhang, Nicholas Roy, and Chuchu Fan. Autotamp: Au-
560 toregressive task and motion planning with llms as translators and checkers. *arXiv preprint*
561 *arXiv:2306.06531*, 2023.
- 562
563 Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong
564 Wu. Seeclck: Harnessing gui grounding for advanced visual gui agents. URL [http://arxiv.org/abs/](http://arxiv.org/abs/2401.10935)
565 [2401.10935](http://arxiv.org/abs/2401.10935).
- 566 Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and
567 Yu Su. Mind2web: Towards a generalist agent for the web, 2023.
- 568
569 Johan Edstedt, Qiyu Sun, Georg Bökman, Mårten Wadenbäck, and Michael Felsberg. RoMa: Robust
570 Dense Feature Matching. 2024.
- 571 Hiroki Furuta, Kuang-Huei Lee, Ofir Nachum, Yutaka Matsuo, Aleksandra Faust, Shixiang Shane
572 Gu, and Izzeddin Gur. Multimodal web navigation with instruction-finetuned foundation models.
573 In *The Twelfth International Conference on Learning Representations*, 2024. URL [https://](https://openreview.net/forum?id=efFmBWioSc)
574 openreview.net/forum?id=efFmBWioSc.
- 575 Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. Exploring network structure, dynamics,
576 and function using networkx. In Gaël Varoquaux, Travis Vaught, and Jarrod Millman (eds.),
577 *Proceedings of the 7th Python in Science Conference*, pp. 11–15.
- 578
579 Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan,
580 and Dong Yu. Webvoyager: Building an end-to-end web agent with large multimodal models,
581 2024.
- 582 Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang,
583 Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. Metagpt: Meta programming for multi-agent
584 collaborative framework. *arXiv preprint arXiv:2308.00352*, 2023.
- 585
586 Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan
587 Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent: A
588 visual language model for gui agents. URL <https://arxiv.org/abs/2312.08914v2>.
- 589
590 Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot
591 planners: Extracting actionable knowledge for embodied agents. In *International Conference on*
592 *Machine Learning*, pp. 9118–9147. PMLR, 2022.
- 593 Hanwen Jiang, Arjun Karapur, Bingyi Cao, Qixing Huang, and Andre Araujo. Omniglue: General-
izable feature matching with foundation model guidance. In *Proceedings of the IEEE/CVF*
Conference on Computer Vision and Pattern Recognition (CVPR), 2024.

- 594 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
595 Narasimhan. Swe-bench: Can language models resolve real-world github issues? URL <http://arxiv.org/abs/2310.06770>.
596
597
- 598 Dongming Jin, Zhi Jin, Xiaohong Chen, and Chunhui Wang. Mare: Multi-agents collaboration
599 framework for requirements engineering, 2024. URL <https://arxiv.org/abs/2405.03256>.
600
- 601 Raghav Kapoor, Yash Parag Butala, Melisa Russak, Jing Yu Koh, Kiran Kamble, Waseem Alshikh,
602 and Ruslan Salakhutdinov. Omniact: A dataset and benchmark for enabling multimodal generalist
603 autonomous agents for desktop and web. URL <http://arxiv.org/abs/2402.17553>.
604
- 605 Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and
606 Ashish Sabharwal. Decomposed prompting: A modular approach for solving complex tasks.
607 In *The Eleventh International Conference on Learning Representations, 2023*. URL https://openreview.net/forum?id=_nGgzQjzaRy.
608
- 609 Avi Kivity, Yaniv Kamay, Dor Laor, Uri Lublin, and Anthony Liguori. kvm: the linux virtual machine
610 monitor. In *Proceedings of the Linux symposium*, volume 1, pp. 225–230. Dttawa, Dntorio, Canada,
611 2007.
- 612 Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham
613 Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating
614 multimodal agents on realistic visual web tasks. URL <http://arxiv.org/abs/2401.13649>.
615
616
- 617 Yihuai Lan, Zhiqiang Hu, Lei Wang, Yang Wang, Deheng Ye, Peilin Zhao, Ee-Peng Lim, Hui Xiong,
618 and Hao Wang. Llm-based agent society investigation: Collaboration and confrontation in avalon
619 gameplay. *arXiv preprint arXiv:2310.14985*, 2023.
- 620 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li,
621 Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer, September 2024a.
622
- 623 Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem.
624 Camel: Communicative agents for " mind" exploration of large scale language model society. 2023.
- 625 Hongxin Li, Jingran Su, Yuntao Chen, Qing Li, and ZHAO-XIANG ZHANG. Sheetcopilot: Bringing
626 software productivity to the next level through large language models. *Advances in Neural*
627 *Information Processing Systems*, 36, 2024b.
628
- 629 Fangru Lin, Emanuele La Malfa, Valentin Hofmann, Elle Michelle Yang, Anthony Cohn, and Janet B.
630 Pierrehumbert. Graph-enhanced large language models in asynchronous plan reasoning. URL
631 <http://arxiv.org/abs/2402.02805>.
- 632 Jiaju Lin, Haoran Zhao, Aochi Zhang, Yiting Wu, Huqiuyue Ping, and Qin Chen. Agentsims: An
633 open-source sandbox for large language model evaluation. *arXiv preprint arXiv:2308.04026*, 2023.
634
- 635 Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng Li, Tianhe Ren, Xueyan Zou, Jianwei Yang,
636 Hang Su, Jun Zhu, Lei Zhang, Jianfeng Gao, and Chunyuan Li. Llava-plus: Learning to use tools
637 for creating multimodal agents, a. URL <https://arxiv.org/abs/2311.05437v1>.
- 638 Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang,
639 Hang Su, Jun Zhu, and Lei Zhang. Grounding dino: Marrying dino with grounded pre-training for
640 open-set object detection, b. URL <https://arxiv.org/abs/2303.05499v4>.
- 641 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,
642 Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui
643 Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang.
644 Agentbench: Evaluating LLMs as agents. In *The Twelfth International Conference on Learning*
645 *Representations, 2024*. URL <https://openreview.net/forum?id=zAdUB0aCTQ>.
646
- 647 Xing Han Lù, Zdeněk Kasner, and Siva Reddy. Weblinx: Real-world website navigation with
multi-turn dialogue. URL <https://arxiv.org/abs/2402.05930v1>.

- 648 Michael M. McKerns, Leif Strand, Tim Sullivan, Alta Fang, and Michael A. G. Aivazis. Building a
649 framework for predictive science. URL <http://arxiv.org/abs/1202.1056>.
- 650
- 651 Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom.
652 Gaia: A benchmark for general ai assistants. URL <http://arxiv.org/abs/2311.12983>.
- 653 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Belle-
654 mare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen,
655 Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra,
656 Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning.
657 518(7540):529–533. ISSN 1476-4687. URL [https://www.nature.com/articles/
658 nature14236](https://www.nature.com/articles/nature14236).
- 659 OpenAI. Gpt-4 omni. <https://openai.com/index/hello-gpt-4o/>, 2024.
- 660
- 661 Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S
662 Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th
663 Annual ACM Symposium on User Interface Software and Technology*, pp. 1–22, 2023.
- 664 Guilherme Potje, Felipe Cadar, Andre Araujo, Renato Martins, and Erickson R Nascimento. Xfeat:
665 Accelerated features for lightweight image matching. In *IEEE/CVF Conference on Computer
666 Vision and Pattern Recognition (CVPR)*, 2024.
- 667
- 668 Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. Android in the
669 wild: A large-scale dataset for android device control, 2023.
- 670
- 671 Christopher Rawles, Sarah Clinckemillie, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Marybeth
672 Fair, Alice Li, William Bishop, Wei Li, Folawiyi Campbell-Ajala, Daniel Toyama, Robert Berry,
673 Divya Tyamagundlu, Timothy Lillicrap, and Oriana Riva. Androidworld: A dynamic benchmarking
674 environment for autonomous agents, 2024. URL <https://arxiv.org/abs/2405.14573>.
- 675 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste
676 Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini
677 1.5: Unlocking multimodal understanding across millions of tokens of context. URL [http:
678 //arxiv.org/abs/2403.05530](http://arxiv.org/abs/2403.05530).
- 679 Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Percy Liang. World of bits:
680 An open-domain platform for web-based agents. In Doina Precup and Yee Whye Teh (eds.),
681 *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings
682 of Machine Learning Research*, pp. 3135–3144. PMLR, 06–11 Aug 2017. URL [https://
683 proceedings.mlr.press/v70/shi17a.html](https://proceedings.mlr.press/v70/shi17a.html).
- 684 Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su.
685 Llm-planner: Few-shot grounded planning for embodied agents with large language models. In
686 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2998–3009, 2023.
- 687
- 688 Liangtai Sun, Xingyu Chen, Lu Chen, Tianle Dai, Zichen Zhu, and Kai Yu. Meta-gui: Towards multi-
689 modal conversational agents on mobile gui. URL <http://arxiv.org/abs/2205.11029>.
- 690 Weihao Tan, Ziluo Ding, Wentao Zhang, Boyu Li, Bohan Zhou, Junpeng Yue, Haochong Xia,
691 Jiechuan Jiang, Longtao Zheng, Xinrun Xu, Yifei Bi, Pengjie Gu, Xinrun Wang, Börje F. Karlsson,
692 Bo An, and Zongqing Lu. Towards general computer control: A multimodal agent for red dead
693 redemption ii as a case study, 2024.
- 694 Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung
695 Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan,
696 Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P. Agapiou, Max Jader-
697 berg, Alexander S. Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David Budden,
698 Yury Sulsky, James Molloy, Tom L. Paine, Caglar Gulcehre, Ziyu Wang, Tobias Pfaff, Yuhuai
699 Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom Schaul,
700 Timothy Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver. Grand-
701 master level in starcraft ii using multi-agent reinforcement learning. 575(7782):350–354. ISSN
1476-4687. URL <https://www.nature.com/articles/s41586-019-1724-z>.

- 702 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlikar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and
703 Anima Anandkumar. Voyager: An open-ended embodied agent with large language models, a.
704 URL <http://arxiv.org/abs/2305.16291>.
705
- 706 Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances
707 large language model capabilities, 2024a. URL <https://arxiv.org/abs/2406.04692>.
- 708 Junyang Wang, Haiyang Xu, Jiabo Ye, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao
709 Sang. Mobile-agent: Autonomous multi-modal mobile device agent with visual perception, b.
710 URL <https://arxiv.org/abs/2401.16158v2>.
- 711 Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai
712 Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents.
713 *Frontiers of Computer Science*, 18(6):1–26, 2024b.
- 714 Luyuan Wang, Yongyu Deng, Yiwei Zha, Guodong Mao, Qinmin Wang, Tianchen Min, Wei Chen,
715 and Shoufa Chen. Mobileagentbench: An efficient and user-friendly benchmark for mobile llm
716 agents, 2024c. URL <https://arxiv.org/abs/2406.08184>.
- 717 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li,
718 Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via
719 multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*, 2023.
- 720 Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe
721 Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou,
722 Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongx-
723 iang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing
724 Huang, and Tao Gui. The rise and potential of large language model based agents: A survey. URL
725 <https://arxiv.org/abs/2309.07864v3>.
726
- 727 Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua,
728 Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio Savarese,
729 Caiming Xiong, Victor Zhong, and Tao Yu. Osworld: Benchmarking multimodal agents for open-
730 ended tasks in real computer environments. URL <http://arxiv.org/abs/2404.07972>.
731
- 732 Mingzhe Xing, Rongkai Zhang, Hui Xue, Qi Chen, Fan Yang, and Zhen Xiao. Understanding
733 the weakness of large language model agents within a complex android environment. URL
734 <http://arxiv.org/abs/2402.06596>.
- 735 Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. Set-of-mark
736 prompting unleashes extraordinary visual grounding in gpt-4v, a. URL <http://arxiv.org/abs/2310.11441>.
737
- 738 John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. Intercode: Standardizing
739 and benchmarking interactive coding with execution feedback, b. URL <http://arxiv.org/abs/2306.14898>.
740
- 741 Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable
742 real-world web interaction with grounded language agents. *Advances in Neural Information*
743 *Processing Systems*, 35:20744–20757, 2022.
744
- 745 Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang, Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qingwei
746 Lin, Saravan Rajmohan, Dongmei Zhang, and Qi Zhang. Ufo: A ui-focused agent for windows os
747 interaction, a. URL <http://arxiv.org/abs/2402.07939>.
748
- 749 Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu.
750 Appagent: Multimodal agents as smartphone users, b. URL <http://arxiv.org/abs/2312.13771>.
751
- 752 Jiwen Zhang, Jihao Wu, Yihua Teng, Minghui Liao, Nuo Xu, Xiao Xiao, Zhongyu Wei, and Duyu
753 Tang. Android in the zoo: Chain-of-action-thought for gui agents, 2024.
754
- 755 Shun Zhang, Zhenfang Chen, Yikang Shen, Mingyu Ding, Joshua B Tenenbaum, and Chuang Gan.
Planning with large language models for code generation. *arXiv preprint arXiv:2303.05510*, 2023.

756 Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,
757 Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic
758 web environment for building autonomous agents. URL [http://arxiv.org/abs/2307.](http://arxiv.org/abs/2307.13854)
759 13854.

760 Kaijie Zhu, Jiaao Chen, Jindong Wang, Neil Zhenqiang Gong, Diyi Yang, and Xing Xie. Dyval:
761 Dynamic evaluation of large language models for reasoning tasks. In *The Twelfth International*
762 *Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=gjfOL9z5Xr)
763 [id=gjfOL9z5Xr](https://openreview.net/forum?id=gjfOL9z5Xr).

764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A BENCHMARK DETAIL

Section A.1 introduces the implementation details and action space settings of the benchmark environments. Section A.2 describes the design logic and implementation of the CRAB framework. Section A.3 describes the our experiment settings in detail. Section A.4 describes the specific format defined in our framework that ease data extension and how to use them. We provides a detailed document to setup experiment environments and reproduce our results.⁴ Fig. 4 shows the structure of modules inside Crab Benchmark-v0.

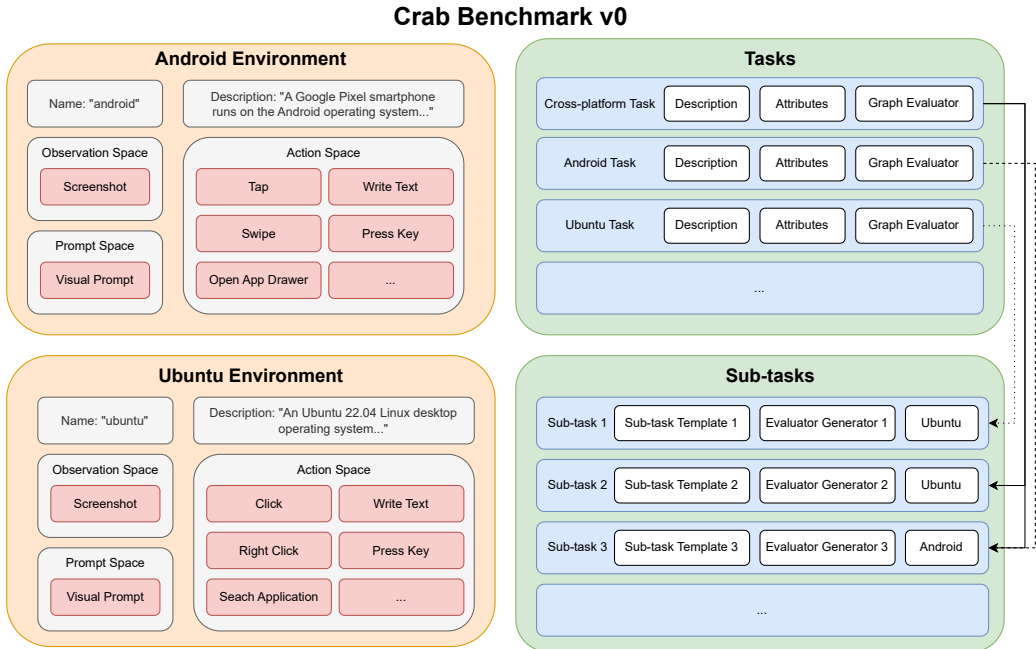


Figure 4: **Module Structure of Crab Benchmark-v0.** The benchmark is divided into two primary sections: the left section, highlighted with warm hues, features two environments, while the right section, accentuated with cool hues, outlines various tasks. Each environment is defined by attributes including name, description, observation space, prompt method, and action space. Blocks marked in red denote actions. As for the tasks, they are composed of multiple sub-tasks and formulated by combine multiple evaluator sub-graphs derived from the sub-task evaluator generators. Arrows illustrate the compositional relationships between tasks and sub-tasks.

A.1 ENVIRONMENT IMPLEMENTATION DETAIL

The Ubuntu environment is launched on a QEMU/KVM (Bellard, 2005; Kivity et al., 2007) Virtual Machine, and the Android environment employs the Google Android Emulator⁵. Interaction with the Ubuntu environment is facilitated using PyAutoGUI⁶ and MSS⁷, which provide high-level commands for mouse and keyboard control and screen capture, respectively. For the Android environment, we use the Android Debug Bridge (ADB)⁸. The detailed action space is described in Table 3.

⁴<https://github.com/camel-ai/crab/blob/main/crab-benchmark-v0/README.md>

⁵<https://developer.android.com/studio/run/emulator>

⁶<https://github.com/asweigart/pyautogui>

⁷<https://github.com/BoBoTiG/python-mss>

⁸<https://developer.android.com/tools/adb>

Table 3: **Action space of Crab Benchmark-v0.** The actions at the top of the table apply to the Ubuntu environment, those in the middle to the Android environment, and those at the bottom are relevant across all environments.

Action Name (Parameters)	Description
<code>click(elem)</code>	Click on <code>elem</code> .
<code>right_click(elem)</code>	Right-click on <code>elem</code> .
<code>double_click(elem)</code>	Double-click on <code>elem</code> .
<code>write_text(text)</code>	Typing the specified text.
<code>press(key)</code>	Press a keyboard key.
<code>hotkey(keys)</code>	Press keyboard keys at the same time.
<code>scroll(direction)</code>	Scrolls page up or down.
<code>search_app(name)</code>	Search for application with name in the system.
<code>tap(elem)</code>	Tap on <code>elem</code> .
<code>long_tap(elem)</code>	Press and hold <code>elem</code> .
<code>swipe(elem, dire, dist)</code>	Swipe from <code>elem</code> in a specified direction and distance.
<code>write_text(text)</code>	Typing the specified text.
<code>press(key)</code>	Press a key, can be <i>home</i> or <i>back</i> .
<code>show_all_drawer()</code>	Show the app drawer to list installed applications.
<code>submit(answer)</code>	Submit <code>answer</code> if needed.
<code>complete()</code>	State that a task is completed.
<code>wait()</code>	Wait the environment to process

A.2 FRAMEWORK DESIGN

CRAB offers a modular and extensible framework for evaluating agent performance in diverse tasks. At the heart of the framework lies the *action*, a unit operation representing the fundamental operation within the benchmark. The *action* is essentially an executable Python function that can be defined with explicit typed parameters and a clear description. *actions* serve not only as building blocks but also as interfaces through which agents interact with the environment. The *evaluator* is a specialized *action* restricted to returning boolean values, signifying the success or failure of an agent’s task. It enhances the *actions* by analyzing the state of the environment and the sequence of *actions* executed by the agent, providing a decisive metric of task accomplishment. Additionally, multiple *evaluators* can be interconnected to form a graph evaluator for complex tasks (Sec. 4.2).

The *benchmark* is a key definition in the framework. A benchmark includes multiple *environments* and cross-environment *tasks*. The *environment* is formed by an action space and an observation space, which are both defined by a list of *actions*, and other essential parameters necessary for its configuration. This composite structure facilitates the execution and monitoring of *actions*, whether on local machines, remote servers, virtual machines, or physical devices networked together. A *task* encapsulates a natural language description and a graph evaluator.

CRAB utilizes Python functions to define all actions and evaluators, embodying a "code as configuration" philosophy. Each function’s docstring outlines its description and parameter definitions, which are then presented to the agent as structured prompts. Compared to traditional methods using data interchange formats like JSON or YAML, Python code configurations provide a more structured approach and fits in modern IDE.

By decoupling actions, environments, tasks, and evaluations, CRAB facilitates a plug-and-play architecture that can adapt to various scenarios. Such a system is scalable, maintainable and expandable, allowing researchers and developers to introduce new tasks and environments without restructuring the entire framework. Our implementation uses *networkx* (Hagberg et al.) for building graph and *dill* (McKerns et al.) for function serialization in our implementation.

A.3 CONFIGURATION BY MODULES

Building on the declarative and modular design of our framework, this section explains the configuration and potential extensibility of each module.

Environment The environments in CRAB are a combination of multiple different uses of actions with some environment metadata, such as name and natural language description. In Crab Benchmark-v0, we use a computer desktop environment and a smartphone environment both based on virtual machine technology. The computer desktop environment, named *Ubuntu*, is installed from an ISO image of Ubuntu 22.04.4 LTS (Jammy Jellyfish) downloaded from the Ubuntu Official website⁹. Necessary applications such as the LibreOffice suite (Writer, Calc, and Impress) and Slack are installed later via snap and apt, according to the task dataset requirements. The smartphone environment, named *Android*, is installed using pre-defined devices (Google Pixel 8 Pro with release name *R*) provided in Google Android Studio¹⁰. We install additional required applications such as *Keep Notes*, *Tasks*, and *Docs* from Google Play. The descriptions of the two environments in Crab Benchmark-v0, which are inserted in the agent prompts, are as follows:

- **Ubuntu:** An Ubuntu 22.04 Linux desktop operating system. The interface displays a current screenshot at each step and primarily supports interaction via mouse and keyboard. You must use searching functionality to open any application in the system. This device includes system-related applications including Terminal, Files, Text Editor, Vim, and Settings. It also features Firefox as the web browser, and the LibreOffice suite—Writer, Calc, and Impress. For communication, Slack is available. The Google account is pre-logged in on Firefox, synchronized with the same account used in the Android environment.
- **Android:** A Google Pixel smartphone runs on the Android operating system. The interface displays a current screenshot at each step and primarily supports interaction through tapping and typing. This device offers a suite of standard applications including Phone, Photos, Camera, Chrome, and Calendar, among others. Access the app drawer to view all installed applications on the device. The Google account is pre-logged in, synchronized with the same account used in the Ubuntu environment.

Action Action implementation in Crab Benchmark-v0 utilize the dynamic feature of Python. It provides an intuitive method to define actions through Python function. Here is an example of action `search_application` in the Ubuntu environment:

```
@action
def search_application(name: str) -> None:
    """Search an application name.

    For exmaple, if you want to open an application named "slack",
    you can call search_application(name="slack"). You MUST use this
    action to search for applications.

    Args:
        name: the application name.
    """
    pyautogui.hotkey("win", "a")
    time.sleep(0.5)
    pyautogui.write(name)
    time.sleep(0.5)
```

Listing 1: Define "search_application" action.

We extract key information from the function through the `@action` decorator as following:

- **Name:** The action name serves as the identifier for backend models. It should semantically match the action’s behavior to improve the accuracy of the agent in executing the action. The function name is extracted as the action name. In this example, `search_application` is the assigned name.
- **Description:** The description provides a natural language explanation of the action to assist the agent in understanding how to use it. The main body of the function’s docstring is used as the description. For example, in this instance, the description outlines the basic usage of the action: *Search an application name*, along with an example of its usage.

⁹<https://releases.ubuntu.com/jammy/ubuntu-22.04.4-desktop-amd64.iso>

¹⁰<https://developer.android.com/studio>

- **Parameters:** The parameters are the arguments that the functions accept, offering flexibility for the agent to control the environment. Typically, a set of parameters is defined, each consisting of a name, type, and a natural language description. Parameters are extracted from the function’s parameters along with their type annotations. Additionally, parameter descriptions are extracted from the `Args` section in the docstring. In this example, there is only one parameter named `name`, with a type of `str`, and its description is the application name.
- **Entry:** The entry represents the implementation of the function, defined within the function body to specify how the action is executed. When the agent invokes the function, the entry is executed with the provided parameters. In this example, we utilize the `pyautogui` package for keyboard control. Initially, it presses a hotkey to enter the application search panel in Ubuntu, then proceeds to type the application name provided by the parameters, finally displaying the search results.

Observation The observation space is represented by a set of actions. These observation actions are designed to be parameter-free and return an observation result. For instance, within the Ubuntu environment, the sole observation action available is the `screenshot` function, defined as follows:

```
@action
def screenshot() -> str:
    """Capture the current screen as a screenshot."""
    with mss() as sct:
        # Capture raw pixels from the screen
        sct_img = sct.grab(sct.monitors[1])
        # Convert to PNG format
        png = tools.to_png(sct_img.rgb, sct_img.size)
        # Encode to Base64 format for easier transmission
        base64_img = base64.b64encode(png).decode("utf-8")
    return base64_img
```

Listing 2: Define the "screenshot" observation action.

This action captures the screen’s current view and encodes it in Base64 format. Additionally, visual prompts are also defined by actions that utilize the output from an observation action as their input, further processing it to generate a visual prompt for the agent.

Evaluator The evaluator in Crab Benchmark-v0 is crafted to assess the outcome of actions performed by the agent within the environment. The evaluator is defined as an action that outputs a boolean value. An example of an evaluator in the Ubuntu environment is the `check_text_in_current_window_name` function, outlined below:

```
@evaluator(env_name="ubuntu")
def check_text_in_current_window_name(text: str) -> bool:
    try:
        out = subprocess.check_output(
            ["xdotool", "getwindowfocus", "getwindowname"], text=True
        ).strip()
    except subprocess.CalledProcessError:
        return False
    return text in out
```

Listing 3: Define "check_text_in_current_window_name" evaluator.

The evaluator function is denoted with an `@evaluator` decorator and specifies its operating environment. The function’s primary role is to execute a check within the system and return a boolean value indicating success or failure based on the condition being evaluated. Here, the function aims to verify whether a specified text appears in the title of the currently focused window. This is achieved through the use of the `subprocess` module to execute system commands that fetch the window’s title, checking if the provided text parameter is contained within it.

Task Following a declarative programming paradigm, the task is defined as a data model. Here is an example of a cross-platform task in the dataset:

```

1026 Task (
1027     id="a3476778-e512-40ca-b1c0-d7aab0c7f18b",
1028     description="Open \{Tasks\} app on Android, check the...",
1029     evaluator=path_graph(
1030         check_current_package_name("com.google.android.apps.tasks"),
1031         check_current_window_process("gnome-control-center"),
1032         check_color_scheme("prefer-dark"),
1033     ),
1034 )

```

Listing 4: Define a task.

In this model, each task is represented as an instance of the `Task` class, which is a subclass of `BaseModel` in *Pydantic*¹¹ package. Each task is uniquely identified by an ID and described by a detailed description. The evaluator component is structured as a graph evaluator, which integrates multiple evaluative functions into a directed graph using the *networkx*¹² package. Each evaluator within this graph must be appropriately parameterized to assess specific conditions relevant to the task. For example, the task demonstrated aims to open the "Tasks" app on Android and perform a series of verifications: it checks whether the correct Android app is opened, whether the current focused window's process name is `gnome-control-center`, and whether the color scheme is set to dark.

Sub-task The sub-task in CRAB is the unit component of in task construction. The following example is a sub-task template that we used to easily generate sub-tasks:

```

1048 SubTask (
1049     id="0f589bf9-9b26-4581-8b78-2961b115ab49",
1050     description="Open \{file_path\} using vim in a terminal, write \{
1051     content\}, then save and exit vim.",
1052     attribute_dict={"file_path": "file_path", "content": "message"},
1053     output_type="file_path",
1054     evaluator_generator=lambda file_path, content: path_graph(
1055         check_current_window_process("gnome-terminal-server"),
1056         is_process_open("vim"),
1057         is_process_close("vim"),
1058         check_file_content(file_path, content),
1059     ),
1060 )

```

Listing 5: Define a task.

In this sub-task model, each sub-task is defined using a similar approach to the main task. The attributes of the sub-task are outlined in an `attribute_dict`, which details the types and roles of each attribute used in the sub-task's operations. The `output_type` field specifies the expected type of output from the sub-task. The types reflected in `attribute_dict` and `output_type`, play a critical role in determining the compatibility and sequential logic of compose multiple sub-tasks. The evaluator for the sub-task is dynamically generated using a lambda function, which crafts an evaluator sub-graph based on the sub-task's attributes.

1068 A.4 TASK DATASET

1069 The task dataset covers a wide range of applications across two platforms, primarily focusing on daily
1070 life, programming, and office work scenarios. The specific applications and software involved are
1071 listed below:
1072

- 1073
- 1074 • Ubuntu: Terminal, Vim, Firefox (Multiple different websites), File Manager, System settings,
- 1075 GIMP, LibreOffice: Impress, Writer, and Calc, VSCode.
- 1076 • Android: Google Map, Keep Notes, Messages, Phone, Google Tasks, Google Calendar,
- 1077 Contacts, Gmail, Google Drive, Files, Settings, Clock, Camera.

1078 ¹¹<https://pydantic.dev/>

1079 ¹²<https://networkx.org/>

We use a JSON format to save the composed tasks, which includes the task ID, overall task description, sub-tasks with their attribute values, and a graph structure represented in an adjacency list. The entire task dataset is defined by the sub-task pool in Python code and the task composition JSON files categorized by task platform.

```

{
  "description": "Combine Image 1 \"/home/crab/Pictures/cat.png\" and
  Image 2 \"/home/crab/assets/campus.png\" using GIMP (GNU Image
  Manipulation Program), placing Image 1 on the left side of Image 2,
  and save the combined image to \"/home/crab/Desktop/background.png\".
  Then, set this combined image as the screen background of the system
  .",
  "tasks": [
    {
      "task": "4cf246ea-0a7f-43da-84b6-61d74a2699af",
      "attribute": {
        "image_path_1": "/home/crab/Pictures/cat.png",
        "image_path_2": "/home/crab/assets/campus.png",
        "output_path": "/home/crab/Desktop/background.png"
      },
      "output": "/home/crab/Desktop/background.png"
    },
    {
      "task": "a207ef38-b3b2-4c6c-a1e3-75c38162f5ba",
      "attribute": {
        "photo_path": "/home/crab/Desktop/background.png"
      },
      "output": null
    }
  ],
  "adjlist": "0 1\n1",
  "id": "d3c917ff-406f-447a-87f5-b8d835cba750"
}

```

Listing 6: Define a composite task in JSON.

B AGENT SYSTEM

B.1 AGENT IMPLEMENTATION

In this section, we outline the implementation of the agents used in our experiments, which leverage advanced multimodal language models from OpenAI, Anthropic, and Google. Each agent is designed to function in multi-environment setups, interacting with various action spaces defined by different environments.

General Framework All agents share a common architecture but are tailored to the specific APIs and capabilities of each language model provider.

Initialization Each agent is initialized with several key parameters, including a description, an action space, the model type, maximum tokens, history message length, and an optional environment description. The initialization process involves:

- **Action Space Conversion:** Actions defined for each environment are converted into a schema compatible with the respective API. This ensures that the actions can be correctly interpreted and executed by the language models.
- **System Message Setup:** Depending on whether the agent is configured for single or multiple environments, a system message is formatted to provide the model with context about the tasks and environments.

Interaction (Chat Method) The core functionality of each agent is encapsulated in its ability to interact with users through a chat method. This involves:

- 1134 • **Content Parsing:** Input content is parsed and formatted to match the requirements of
1135 the respective API. This includes structuring user messages and any necessary contextual
1136 information.
- 1137 • **Request Construction:** The request payload is constructed, incorporating the system
1138 message, chat history, and the newly parsed user input.
- 1139 • **API Communication:** The constructed request is sent to the appropriate API, which
1140 generates a response. The agents handle API-specific constraints such as rate limits and
1141 response formats.
- 1142 • **Response Handling:** The response from the API is processed to extract any tool calls
1143 suggested by the model. These are then appended to the chat history, maintaining a coherent
1144 conversation state.

1146 **Multi-Environment Support** For agents configured to operate in multiple environments, additional
1147 logic ensures that actions are correctly associated with their respective environments. This involves
1148 modifying action names and descriptions to reflect their environmental context and handling responses
1149 accordingly.

1151 **Utilities and Shared Functions** Several utility functions support the operation of these agents, facil-
1152 itating tasks such as content parsing, action prompt generation, and schema conversion. These shared
1153 functions ensure consistency and reduce redundancy across the different agent implementations.

1155 B.2 INTER-AGENT COMMUNICATION STRATEGIES

1156 In this section we introduce the details of two multi-agent communications methods, which are
1157 introduced in 6.1.

1159 **Multi-agent Communication by Functionality** This setting involves two agents: a main agent
1160 prompted with the task description and a tool agent with the entire action space. The main agent
1161 generates the instruction for the next step and sends it to the tool agent. The tool agent chooses the
1162 proper action with parameters and a target environment, then feeds it back to the system.

1164 **Multi-agent Communication by Environment** This setting involves four agents in our benchmark
1165 setting: a main agent prompted with the task description and three tool agents, each corresponding to
1166 the environments of Android, Ubuntu, and Root, with the respective action spaces. The main agent
1167 generates the instruction for the next step and sends it to the tool agents. Each sub-environment
1168 agent receives the message containing the instruction and environment observation information.
1169 The environment agents process the message using their specialized models and action schemas,
1170 performing the required actions within their environments.

1171 B.3 AGENT PROMPT

1173 B.3.1 SINGLE AGENT

1175 Prompt

1176 You are a helpful assistant. Now you have to do a task as described below:
1177 `**{task_description}**`.
1178 You should never forget this task and always perform actions to achieve this task. And this is
1179 the description of each given environment: `{env_description}`. A unit operation you
1180 can perform is called action in a given environment. For each environment, you are given a
1181 limited action space as function calls:
1182 `{action_descriptions}`
1183 You may receive a screenshot of the current system. You may receive a screenshot of a
1184 smartphone app. The interactive UI elements on the screenshot are labeled with numeric tags
1185 starting from 1.
1186 In each step, You MUST explain what do you see from the current observation and the plan of
1187 the next action, then use a provided action in each step to achieve the task. You should state

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

what action to take and what the parameters should be. Your answer **MUST** be a least one function call. You **SHOULD NEVER** ask me to do anything for you. Always do them by yourself using function calls.

Prompt

You are a helpful assistant. Now you have to do a task as described below:
`**{task_description}**`
 You should never forget this task and always perform actions to achieve this task. And this is the description of each given environment: `{env_description}`. You will receive screenshots of the environments. The interactive UI elements on the screenshot are labeled with numeric tags starting from 1.
 A unit operation you can perform is called Action. You have a limited action space as function calls: `{action_descriptions}`. You should generate JSON code blocks to execute the actions. Each code block **MUST** contains only one json object, i.e. one action. You can output multiple code blocks to execute multiple actions in a single step. You must follow the JSON format below to output the action.
`{"name": "action_name", "arguments": {"arg1": "value1", "arg2": "value2"}}`
 or if not arguments needed:
`{"name": "action_name", "arguments": {}}`
 You **MUST** use exactly the same "action_name" as I gave to you in the action space. You **SHOULDN'T** add any comments in the code blocks.
 In each step, You **MUST** explain what do you see from the current observation and the plan of the next action, then use a provided action in each step to achieve the task. You should state what action to take and what the parameters should be. Your answer **MUST** contain at least one code block. You **SHOULD NEVER** ask me to do anything for you. Always do them by yourself.

B.3.2 MULTI-AGENT BY FUNCTIONALITY

Main Agent Prompt

You are a helpful assistant. Now you have to do a task as described below: `{task_description}`. And this is the description of each given environment: `{env_description}`. A unit operation you can perform is called action in a given environment. For each environment, you are given a limited action space as function calls: `{action_descriptions}`
 You may receive a screenshot of the current system. The interactive UI elements on the screenshot are labeled with numeric tags starting from 1. For each step, You must state what actions to take, what the parameters are, and you **MUST** provide in which environment to perform these actions.

Tool Agent Prompt

You are a helpful assistant in generating function calls. I will give you a detailed description of what actions to take next, you should translate it into function calls. please do not output any other information.

1242 B.3.3 MULTI-AGENT BY ENVIRONMENT

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

Main Agent Prompt

You are a main agent, and your goal is to plan and give instructions to sub-agents in each environment to complete the final task. Now you have to do a task as described below: {description}. The description of each given environment: {env_description}. For each step, you are required to provide high-level instructions detailing the next actions to be taken. Additionally, you must specify which sub-agent in the designated environment should execute these instructions. If a sub-agent is not needed for a particular step, you may instruct it to skip that step.

Root Environment Agent Prompt

You are a sub-agent responsible for the crab benchmark root environment. Your goal is to assist the main agent in completing the whole task: "{description}". You can only complete the task or submit the result when the main agent tells you the whole task has been completed. Otherwise, you can only call SKIP.

Sub-environment Agent Prompt

You are a sub-agent responsible for the {environment} environment. The description of the {environment} environment is: {env_description}. Your goal is to assist the main agent in completing the final task by performing actions in the {environment} environment according to the instructions from the main agent. The final task is described below: {task_description}. A unit operation you can perform is called action in a given environment. You can only execute action in the {environment} environment. For the {environment} environment, you are given a limited action space as function calls: {action_descriptions}. The interactive UI elements on the screenshot are labeled with numeric tags starting from 1. For each step, You will receive an instruction telling you what you need to do next. After analyzing the instruction you received and the current {environment} system, if you think you don't need to do anything in the current {environment} system, you should choose SKIP action. Otherwise, you must state what actions to take, what the parameters are, and you MUST provide in which environment to perform these actions. Your answer must be function calls. Please do not output any other information. You must make sure all function calls get their required parameters.

C FURTHER RESULT ANALYSIS

This section further discusses our experimental results in detail. Section C.1 categorizes the results into three types of tasks: Ubuntu, Android, and cross-platform, and provides further analysis. Section C.3 examines three specific tasks and analyzes the performance of different agent settings on each.

C.1 RESULT BY PLATFORMS

Table 4, 5 and 6 show the experiment results on Ubuntu Tasks, Android Tasks, and cross-platform Tasks, respectively.

We find that certain models demonstrate a distinct preference or better alignment with specific platforms. The GPT-4o, Gemini, and Claude models, for instance, show notably better outcomes on Android platforms. This suggests potential optimizations or intrinsic features within these models that cater effectively to the Android environment's requirements. Conversely, the GPT-4 Turbo model exhibits superior performance on Ubuntu tasks, hinting at possible architectural or training aspects that are better suited for that specific environment.

Cross-platform tasks necessitate functionality across different operating systems or platforms, demand a broader capability range and more sophisticated agent coordination. The importance of CR is

Table 4: Evaluation results on Ubuntu tasks.

Agent system		Metrics				Termination Reason		
Model	Structure	SR(%) ↑	CR(%) ↑	EE(%) ↑	CE(%) ↑	FC(%)	RSL(%)	IA(%)
GPT-4o	Single	9.59	30.82	3.22	4.87×10^{-4}	6.85	58.90	24.66
GPT-4o	By Func	9.59	24.20	2.72	4.30×10^{-4}	5.48	63.01	21.92
GPT-4o	By Env	10.96	22.88	2.74	2.29×10^{-4}	5.48	43.84	39.73
GPT-4 TURBO	Single	10.96	31.09	4.08	5.57×10^{-4}	2.74	65.75	20.55
GPT-4 TURBO	By Func	12.33	28.95	3.70	4.18×10^{-4}	8.22	32.88	46.58
GEMINI 1.5 PRO	Single	1.37	7.76	0.63	n/a	0.00	47.95	50.68
GEMINI 1.5 PRO	By Func	1.37	3.31	0.33	n/a	0.00	20.55	78.08
CLAUDE 3 OPUS	Single	0.00	9.54	0.72	0.63×10^{-4}	8.22	58.90	32.88
CLAUDE 3 OPUS	By Func	0.00	4.93	0.46	0.47×10^{-4}	27.40	34.25	38.36
GPT-4o w/o FC	Single	10.96	22.58	2.30	4.49×10^{-4}	5.48	54.79	28.77
PIXTRAL-12B	Single	0.00	2.97	0.22	0.24×10^{-4}	1.37	80.82	17.81
LLAVA-OV-72B	Single	0.00	3.31	0.20	0.35×10^{-4}	17.81	64.38	17.81

especially critical in such environments, where it serves as a more reliable metric for distinguishing between agent models than SR. Given the presence of all Gemini, Claude, and open source model agents’ SR is 0.0, indicating that Completion Ratio more effectively captures an agent model’s capability, thereby better reflecting its robustness and adaptability to complex requirements. On cross-platform tasks, GPT-4 Turbo (Single) exhibits a CR of 52.61%, which indicates that even though SR might be lower, the agent covers a significant portion of task objectives before termination.

Furthermore, analyzing the reasons for task termination offers additional insights into the operational challenges these models encounter. False Completion is notably prevalent in Android tasks. Reach Step Limit remains the most frequent cause of termination, particularly in cross-platform tasks. The Claude model exhibits a significantly high Invalid Action ratio in cross-platform tasks, indicating its difficulties in managing multi-environment scenarios effectively. The GPT-4o with JSON mode shows a extremely high IA ratio in Android tasks, proving the serious hallucination problem under this setting.

Overall, these findings underscore the necessity of selecting the appropriate agent model and configuration based on specific platform and task needs. The variability in model performance across different setups also highlights the ongoing need for development and refinement of multi-agent systems to enhance their versatility and efficacy in increasingly diverse and complex operational environments. These results comparing SR and CR also demonstrates the important of our graph evaluator in agent evaluation.

C.2 COMPARISON BETWEEN SINGLE AGENT AND MULTI-AGENT

The experimental results indicate that multi-agent structures perform slightly worse than single-agent systems, which is somewhat unusual. We analyse the possible reasons here.

First, comparing in False Completion Rate, we attribute the lower Success Rate (SR) of Multi-agent to a high False Completion Rate—where the agent incorrectly assumes that the task is complete. As observed in failure cases (e.g., the Cross-platform Task case study in Appendix C.3), Sub-agents often misinterpret the Main agent’s instructions. Despite being required to perform a final action, the instructions lead Sub-agents to prematurely conclude that the task is complete, resulting in incorrect “complete” actions. While this issue also occurs in Multi-Env, it happens less frequently. By analysing the communication logs, we believe this is due to information loss during inter-agent communication. Sometimes, the main agent gives a correct instruction, but the sub-agent misunderstands it because it does not have the context. Natural language, while effective for aligning with human understanding in LLM communication, is less suited for inter-agent communication, leading to information loss during compression and interpretation, which weakens the performance of multi-agent structures.

Next, comparing in Invalid Action Rate, we observe that in single-platform tasks, both Multi-Env and Multi-Func suffer from similar inter-agent communication issues, as indicated by their high Invalid

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

Table 5: Evaluation results on Android tasks.

Agent system		Metrics				Termination Reason		
Model	Structure	SR(%) ↑	CR(%) ↑	EE(%) ↑	CE(%) ↑	FC(%)	RSL(%)	IA(%)
GPT-4o	Single	24.14	47.91	5.84	7.17×10^{-4}	13.79	58.62	3.45
GPT-4o	By Func	24.14	48.74	6.83	9.19×10^{-4}	24.14	37.93	13.79
GPT-4o	By Env	27.59	53.34	6.99	4.58×10^{-4}	13.79	44.83	13.79
GPT-4 TURBO	Single	6.90	27.08	2.60	2.87×10^{-4}	20.69	62.07	10.34
GPT-4 TURBO	By Func	20.69	37.01	5.00	5.92×10^{-4}	13.79	51.72	13.79
GEMINI 1.5 PRO	Single	17.24	34.52	4.82	n/a	10.34	65.52	6.90
GEMINI 1.5 PRO	By Func	17.24	35.99	4.31	n/a	31.03	37.93	13.79
CLAUDE 3 OPUS	Single	13.79	41.90	5.07	5.37×10^{-4}	20.69	55.17	10.34
CLAUDE 3 OPUS	By Func	13.79	44.02	4.75	5.35×10^{-4}	48.28	31.03	6.90
GPT-4o w/o FC	Single	10.34	14.29	1.72	2.94×10^{-4}	3.45	6.90	79.31
PIXTRAL-12B	Single	3.45	24.17	2.16	2.72×10^{-4}	0.00	65.52	31.03
LLAVA-OV-72B	Single	3.45	13.51	1.36	3.00×10^{-4}	3.45	93.10	0.00

Table 6: Evaluation results on cross-platform tasks.

Agent system		Metrics				Termination Reason		
Model	Structure	SR(%) ↑	CR(%) ↑	EE(%) ↑	CE(%) ↑	FC(%)	RSL(%)	IA(%)
GPT-4o	Single	16.67	51.24	5.21	3.98×10^{-4}	5.56	38.89	38.89
GPT-4o	By Func	22.22	50.00	4.15	3.13×10^{-4}	11.11	44.44	22.22
GPT-4o	By Env	5.56	43.54	3.22	1.60×10^{-4}	11.11	72.22	11.11
GPT-4 TURBO	Single	5.56	52.61	4.60	2.89×10^{-4}	11.11	66.67	16.67
GPT-4 TURBO	By Func	5.56	46.17	4.06	2.67×10^{-4}	16.67	50.00	27.78
GEMINI 1.5 PRO	Single	0.00	16.14	1.15	n/a	0.00	72.22	27.78
GEMINI 1.5 PRO	By Func	0.00	13.65	1.21	n/a	5.56	77.78	16.67
CLAUDE 3 OPUS	Single	0.00	24.50	1.93	1.24×10^{-4}	0.00	55.56	44.44
CLAUDE 3 OPUS	By Func	0.00	18.96	1.93	1.20×10^{-4}	0.00	38.89	61.11
GPT-4o w/o FC	Single	0.00	39.11	3.51	3.28×10^{-4}	5.56	50.00	44.44
PIXTRAL-12B	Single	0.00	12.35	0.62	0.44×10^{-4}	0.00	72.22	27.78
LLAVA-OV-72B	Single	0.00	9.07	0.48	0.53×10^{-4}	5.56	66.67	27.78

Action rates. However, in cross-platform tasks (Table 6), the Single agent’s Invalid Action rate is significantly higher than that of the Multi-agent by environment structures on GPT-4o model. Cross-platform tasks require frequent environment changes with varying action spaces, and if the model’s performance output is inadequate, it often generates correct actions in the wrong environment, invalid actions in the correct environment, or correct actions in correct environment but in the wrong format. This phenomenon highlights the limitations of current general-purpose LLMs, where multi-agent structures can be advantageous. By assigning each agent a specific responsibility and a limited action space, multi-agent structures can mitigate these issues.

To improve multi-agent system performance, we suggest to follow two approaches: (1) Developing better multi-agent structures to minimize information loss during communication, and (2) Introducing a critical agent to correct hallucinations or information loss during communication. These improvements, however, come with a trade-off, namely an increase in token costs within the agent system. Within our benchmark framework, users can utilize the error log we provide to analyze the bottlenecks of their agents and refine their designs.

C.3 CASE STUDY

To better understand how different agents perform the same task and exhibit varied properties, we present visual results along with detailed metrics and logs for three cases by platform. The screenshots illustrate the progress of agents executing tasks according to specific natural language instructions.

1404 C.3.1 CROSS-PLATFORM TASK
1405

1406 **Task: Open the "Tasks" app on an Android device, check the first incomplete task, and then**
 1407 **execute it as described.** The first task, found incomplete in the "Tasks" app, involves **switching the**
 1408 **system to dark mode in Ubuntu via the "Settings" application.**

1409 This task exemplifies message passing across different environments, where the "incomplete task"
 1410 serves as the critical information that the agent must relay and apply in the Ubuntu setting. These
 1411 two phases—retrieving the task details via the phone and executing the task on a computer—are
 1412 inseparably linked and cannot be treated as distinct tasks. The agent can only proceed to the second
 1413 stage after successfully acquiring information from the first.

1414 In this task, GPT-4o (single agent), GPT-4 Turbo (single agent), and GPT-4 Turbo (multi-agent by
 1415 functionality) all successfully complete the task using the minimal steps necessary to locate and exe-
 1416 cute the task, demonstrating their efficiency in managing multiple environments simultaneously. On
 1417 the other hand, both GPT-4o (multi-agent by functionality) and GPT-4o (multi-agent by environment)
 1418 also perform commendably, completing the task up until the final step. However, after incorrectly
 1419 performing the last step, they both erroneously conclude the task is completed and exit. This indicates
 1420 a communication breakdown, where the sub-agents misinterpret the instructions from the main agent.
 1421 The remaining four agents fail to complete the task. Agents equipped with the Gemini model do
 1422 not even manage to open the "Tasks" app within the allocated step limit, whereas agents with the
 1423 Claude model quickly open the "Tasks" app to complete the first step but fail at the task execution.
 1424 The performance disparity between single-agent and multi-agent configurations in both the Gemini
 1425 and Claude models highlights the variance in capability across different models and devices.

1426 C.3.2 UBUNTU TASK
1427

1428 **Task: Create a new directory "/home/crab/assets_copy" and copy all files with the specified**
 1429 **"txt" extension from "/home/crab/assets" to the directory "/home/crab/assets_copy".**

1431 This task can be approached through multiple methods. An agent may opt for a straightforward
 1432 strategy first using the `search_application` command to find the Terminal, then using Linux
 1433 commands to create the directory and copy the necessary files. Alternatively, the agent could employ
 1434 a GUI-based approach, manually creating the folder and selecting files through actions like `click`
 1435 and `right_click`. We evaluate various agent systems in a single-agent setting for this task. As
 1436 illustrated in Table 7–10, both GPT-4o and GPT-4 Turbo from OpenAI successfully interpret the task
 1437 instructions and employ a simpler solution using Terminal commands. These agents also demonstrate
 1438 superior capability in understanding the UI, selecting the correct commands, and accurately using the
 1439 Terminal application to fulfill the task requirements.

1440 Conversely, the Gemini and Claude agents, despite attempting to solve the task with Terminal,
 1441 ultimately fail in different ways. Both agents struggle with precise clicking and selecting the correct
 1442 icons for the intended actions, even though they share the same visual prompting mechanism as
 1443 GPT-4o and GPT-4 Turbo. For instance, the Claude agent mistakenly opens the Ubuntu Desktop
 1444 Guide instead of the Terminal and continues executing commands in the wrong application without
 1445 realizing the error. The Gemini agent, on the other hand, unexpectedly opens the Firefox browser
 1446 before correctly navigating to the Terminal but still interacts incorrectly with unrelated applications
 1447 and icons. Unlike Claude, Gemini does not type in commands in the wrong applications but persists
 1448 in exploring alternative methods using the Files application’s UI. Despite taking significantly more
 1449 steps than the GPT-4o and GPT-4 Turbo agents, neither the Claude nor the Gemini agents achieve the
 1450 task’s goal.

1451 C.3.3 ANDROID TASK
1452

1453 **Task: In Android, using the "Contacts" app, find the email of the contact named John Lauphin,**
 1454 **then using the "Gmail" app, send an email to that contact with the subject "Hello John."**

1455 This task consists of sub-tasks across two different applications. Agents must sequentially open the
 1456 two apps, retrieve the email address from the first app, and use it in the second app to send an email.
 1457 This straightforward yet formal task can be completed using various methods. Agents may need to
 locate the contact in the Contacts app and then use the retrieved email address to send a message. We

1458 reports the performance of agents in a multi-agent setting for this challenging task. Following is the
1459 details of agents in operating the task.
1460

1461 **GPT-4o multi-agent by functionality** In steps 1-11, the agent tries to open the Contacts app but
1462 mistakenly opens Google Assistant multiple times. In steps 12-14, the agent successfully enters the
1463 Contacts app and finds the contact information. The agent then returns to the home page, and the
1464 process is terminated due to the limitation of operation steps.
1465

1466 **GPT-4 Turbo multi-agent by functionality** In steps 1-2, the agent tries to open the Contacts app
1467 but mistakenly opens Google Messages. In steps 3-5, the agent opens the Contacts app and obtains the
1468 corresponding information. In steps 6-14, the agent repeatedly opens Google Chrome and Messages
1469 apps, failing to find the Gmail app as planned.
1470

1471 **Gemini 1.5 Pro multi-agent by functionality** In steps 1-2, the agent finds the Contacts app and
1472 enters it. However, the agent misunderstands the instruction, gets lost in creating a new contact with
1473 the given name, and cannot obtain the corresponding information.
1474

1475 **Claude 3 Opus multi-agent by functionality** In steps 1-7, the agent tries to open the Contacts
1476 app but mistakenly opens Google Messages multiple times. In steps 7-11, the agent tries to open the
1477 Contacts app but mistakenly opens Google Assistant. In steps 12-14, the agent successfully enters
1478 the Contacts app and finds the contact information. The agent then returns to the home page, plans to
1479 open the Gmail app, and the process is terminated due to the limitation of operation steps.
1480

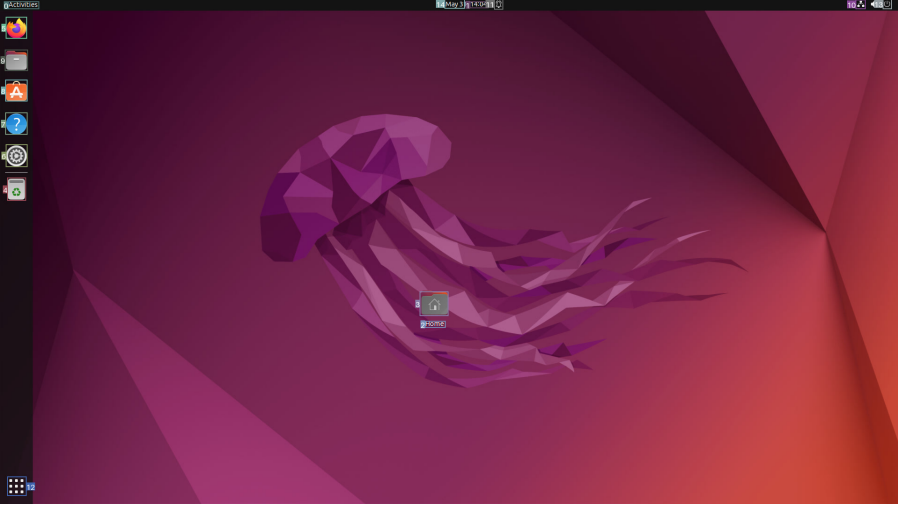
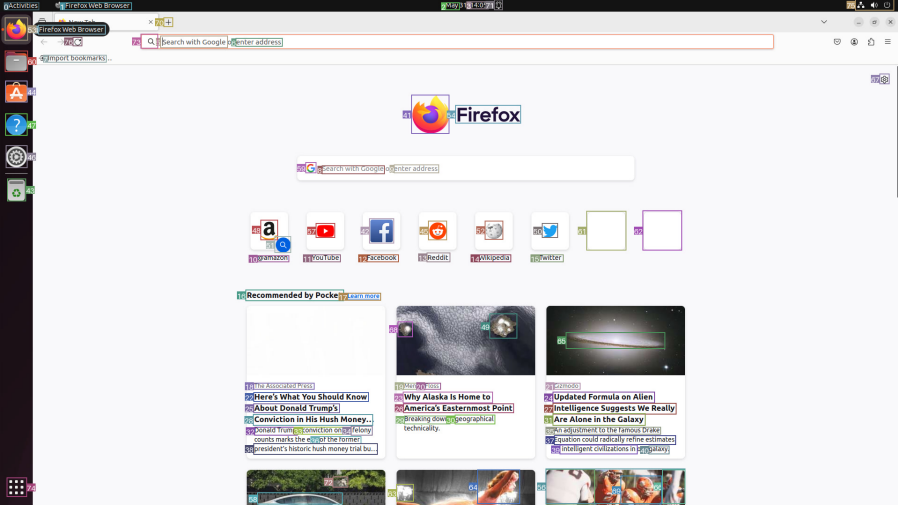
1481 **GPT-4o multi-agent by environment** In steps 1-7, the agent plans to open the Contacts app, but
1482 the operation fails due to an error in opening the app drawer, which prevents the agent from finding
1483 and tapping the Contacts app. In steps 8-11, the agent successfully enters the Contacts app and
1484 obtains the information. In steps 12-14, the agent opens the Gmail app, navigates to the sending page,
1485 and tries to input the retrieved email address as the recipient.
1486

1487 **Analysis** For the agents which are organized by functionality, Gemini 1.5 Pro struggles to complete
1488 the first operation. Although it recognizes and opens the Contacts app as instructed, it fails to proceed
1489 further. In contrast, Claude 3 Opus and GPT-4o successfully obtain the necessary information. In
1490 the initial phase, the multiple agents agree that opening the Contacts app is the first step. However,
1491 they often fail to find the correct position to tap, frequently opening incorrect apps such as Google
1492 Assistant and Messages. Once the agents do open the correct app, they usually find the email address
1493 of the contact quickly. Even when agents plan to go back home and open the Gmail app to send the
1494 message, due to the limitation of operations, the system ended. As shown in steps 3-5, GPT-4 Turbo
1495 quickly finishes the corresponding task after opening the correct app. However, similar to GPT-4o,
1496 GPT4-Turbo agents get stuck as they can not open the correct apps in the following steps. Besides,
1497 GPT-4o (multi-agent by environment) overcomes the issue encountered by GPT-4o (multi-agent by
1498 functionality). Even affected by not being able to access the app drawer, the system could still find
1499 and copy the corresponding information and change to the Gmail app for further operations.
1500

1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565


Table 7: **Ubuntu task case with Gemini (Single):** Create a new directory "/home/crab/assets_copy" and copy all files with the specified "txt" extension from "/home/crab/assets" to the directory "/home/crab/assets_copy".

Step	Agent Observation and Action
0	 <p data-bbox="802 898 922 926">click (5)</p>
1	 <p data-bbox="802 1451 922 1478">click (1)</p>

Continued on next page

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

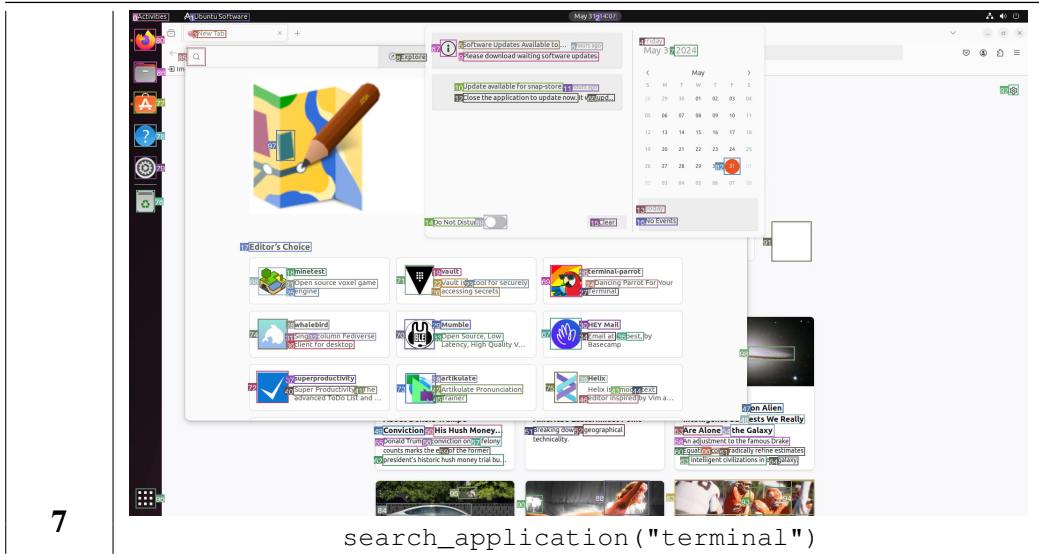
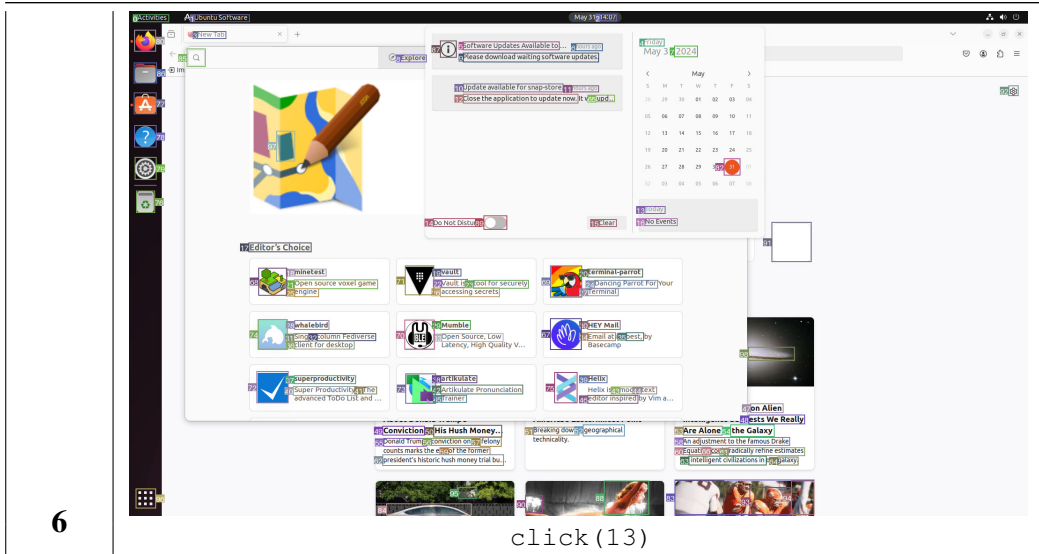
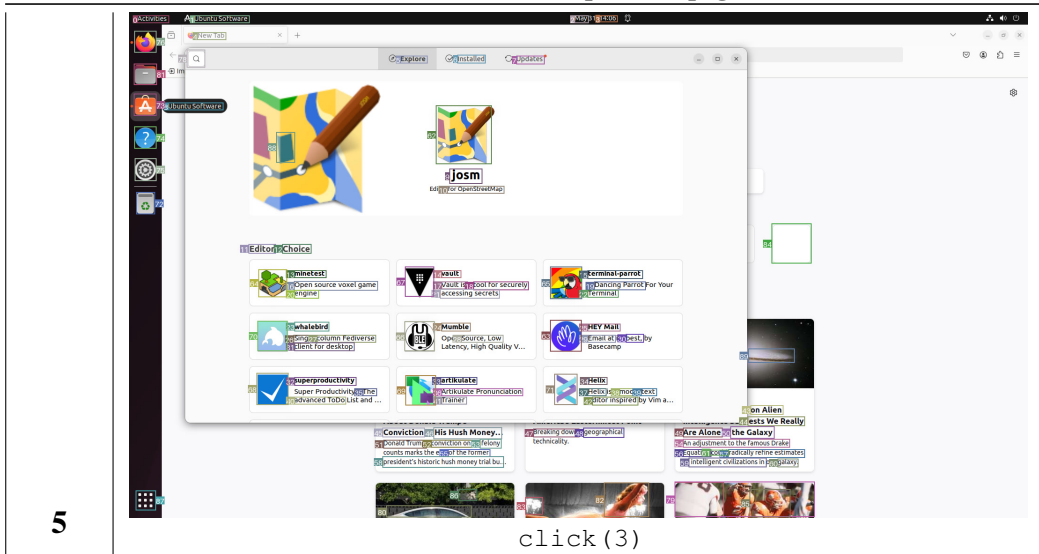
Table 7 continued from previous page

2	 A screenshot of the Firefox browser home page. The address bar is empty. Below the address bar are social media icons for Amazon, YouTube, Facebook, Reddit, Wikipedia, and Twitter. A 'Recommended by Pocket' section is visible, featuring several article thumbnails with titles like 'Here's What You Should Know About Donald Trump's Conviction in His Hush Money...', 'Why Alaska is Home to America's Eastermost Point', and 'Updated Formula on Alien Intelligence Suggests We Really Are Alone in the Galaxy'.
	click (4)
3	 A screenshot of the Firefox browser home page, identical to the previous one. The search bar in the address bar now contains the text 'search with google or enter address'.
	search_application ("terminal")
4	 A screenshot of a terminal window. The terminal has a dark background and shows a prompt character '>' followed by a cursor. Below the terminal, there are several window titles for 'Gnome Terminal Mark 1' and 'Gnome Terminal Mark 2'.
	click (9)

Continued on next page

Table 7 continued from previous page

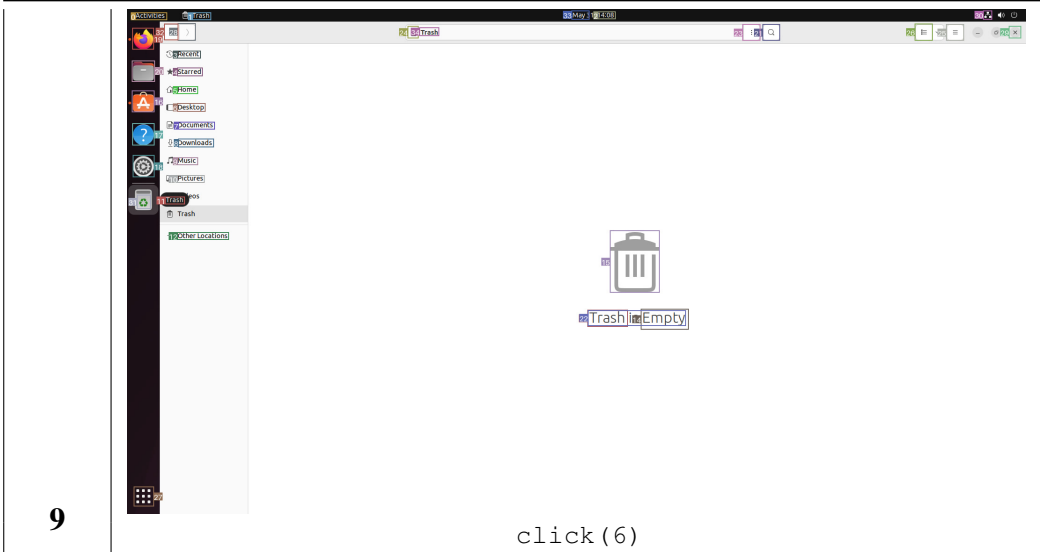
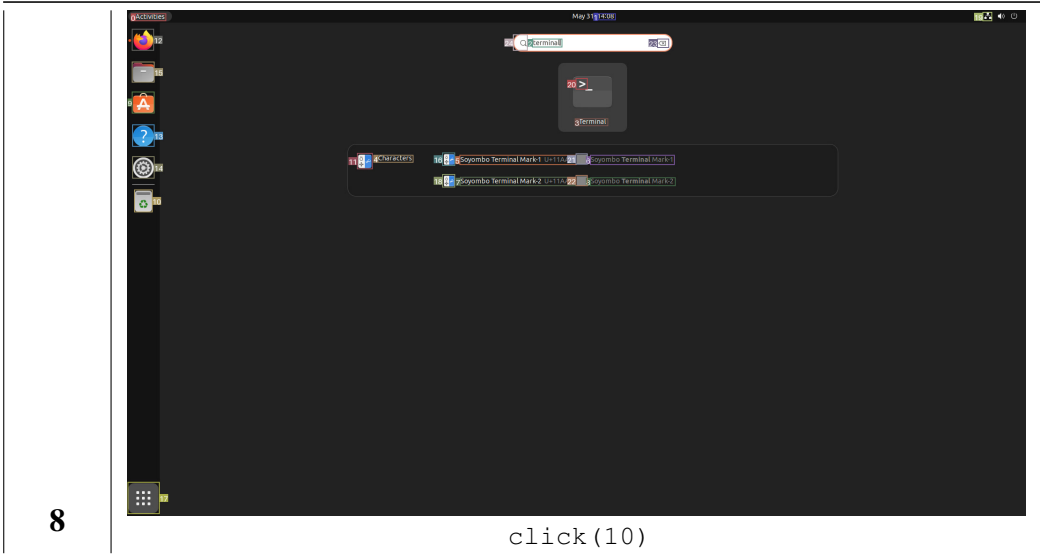
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673



Continued on next page

Table 7 continued from previous page

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

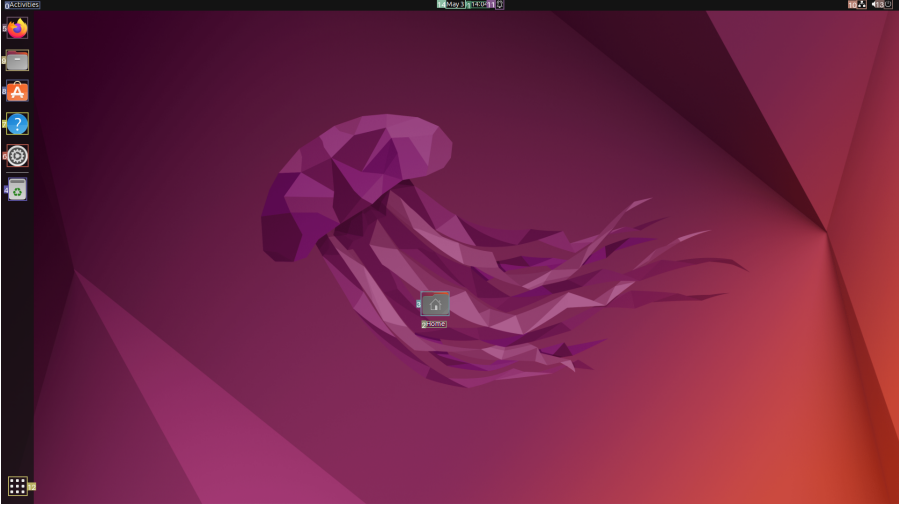
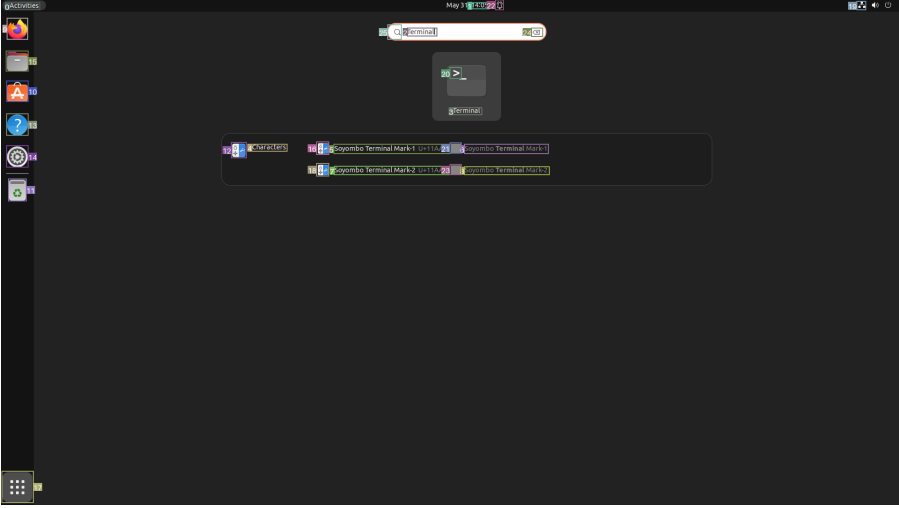


10-14 The agent is stuck at this stage and keeps clicking useless elements.

Completed Nodes / Total Nodes: 0 / 2 Termination Reason: Reach Step Limit

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

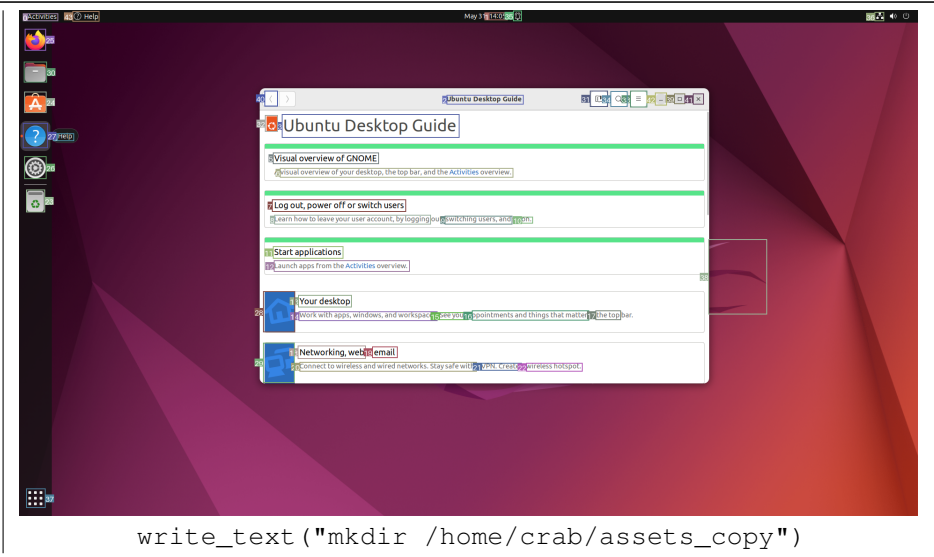
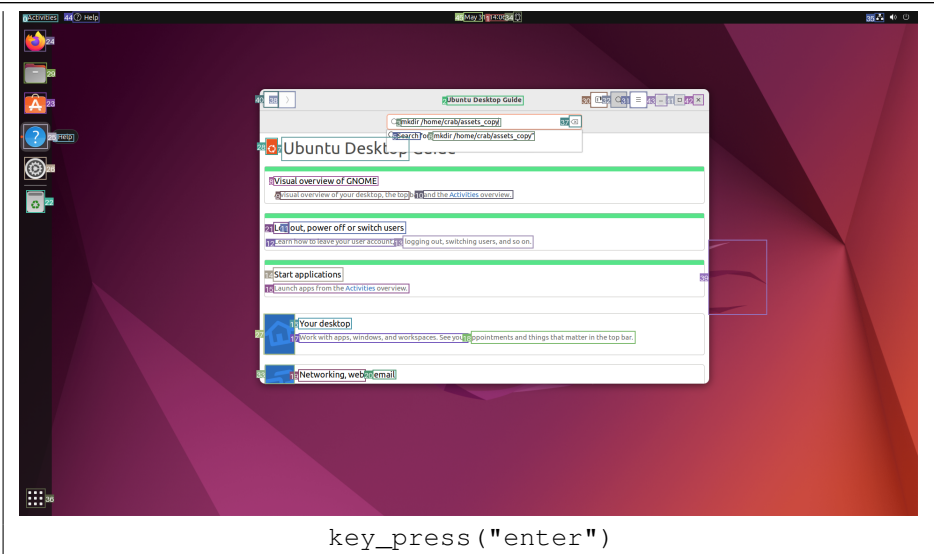
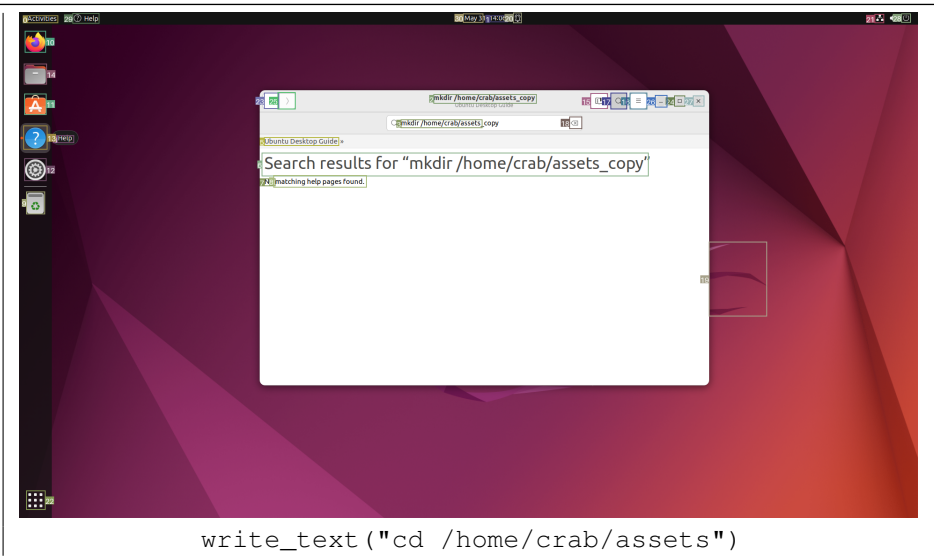
Table 8: **Ubuntu task case with Claude (Single):** Create a new directory "/home/crab/assets_copy" and copy all files with the specified "txt" extension from "/home/crab/assets" to the directory "/home/crab/assets_copy".

Step	Agent Observation and Action
<p>0</p>	 <p>search_application("terminal")</p>
<p>1</p>	 <p>click(13)</p>

Continued on next page

Table 8 continued from previous page

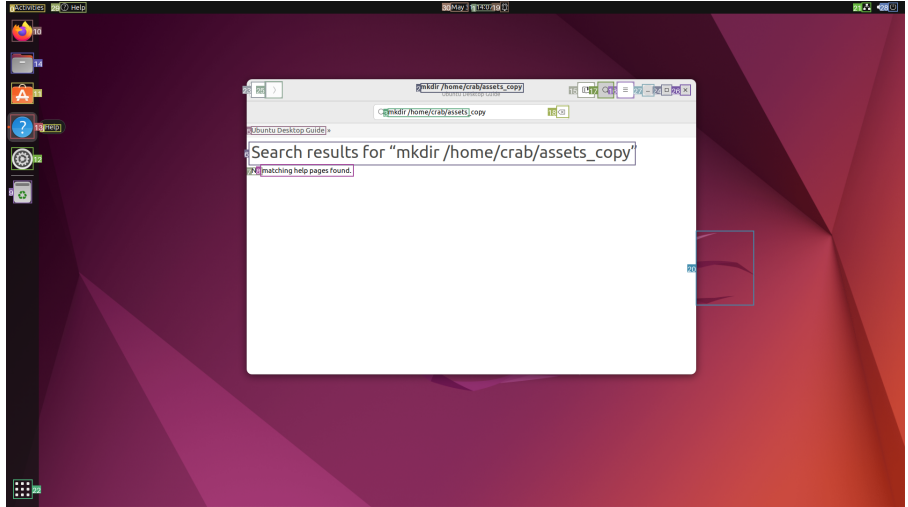
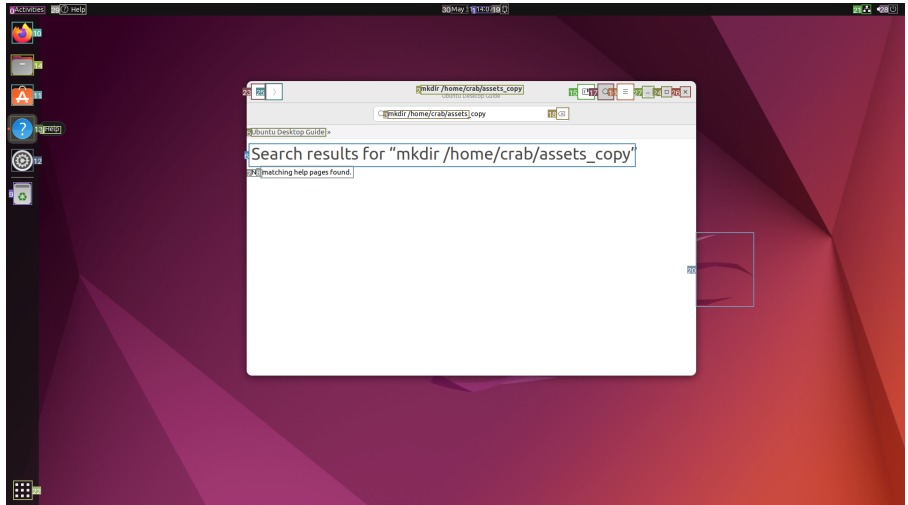
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835

2	 A screenshot of an Ubuntu desktop environment. A window titled "Ubuntu Desktop Guide" is open, displaying sections like "Visual overview of GNOME", "Log out, power off or switch users", "Start applications", "Your desktop", and "Networking, web, email". <pre>write_text("mkdir /home/crab/assets_copy")</pre>
3	 A screenshot of the same Ubuntu desktop environment. A terminal window is open, showing the command <code>mkdir /home/crab/assets_copy</code> being entered and executed. <pre>key_press("enter")</pre>
4	 A screenshot of the same Ubuntu desktop environment. A search results window is open, displaying "Search results for 'mkdir /home/crab/assets_copy'" and "2 matching help pages found". <pre>write_text("cd /home/crab/assets")</pre>

Continued on next page

Table 8 continued from previous page

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

5	 <p>Search results for "mkdir /home/crab/assets_copy" 770 [matching help pages found]</p>
6	 <p>Search results for "mkdir /home/crab/assets_copy" 770 [matching help pages found]</p>

`key_press("enter")`



`write_text("mkdir assets_copy")`

7-14 | The agent is stuck at this stage and keeps pressing keys.

Completed Nodes / Total Nodes: 0 / 2 Termination Reason: Reach Step Limit

1890
 1891
 1892
 1893
 1894
 1895
 1896
 1897
 1898
 1899
 1900
 1901
 1902
 1903
 1904
 1905
 1906
 1907
 1908
 1909
 1910
 1911
 1912
 1913
 1914
 1915
 1916
 1917
 1918
 1919
 1920
 1921
 1922
 1923
 1924
 1925
 1926
 1927
 1928
 1929
 1930
 1931
 1932
 1933
 1934
 1935
 1936
 1937
 1938
 1939
 1940
 1941
 1942
 1943

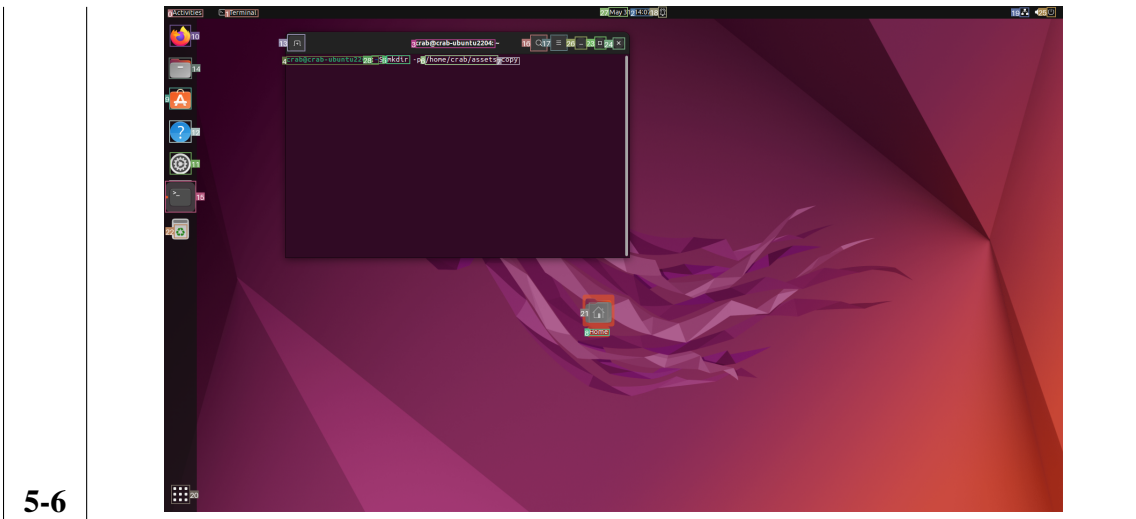
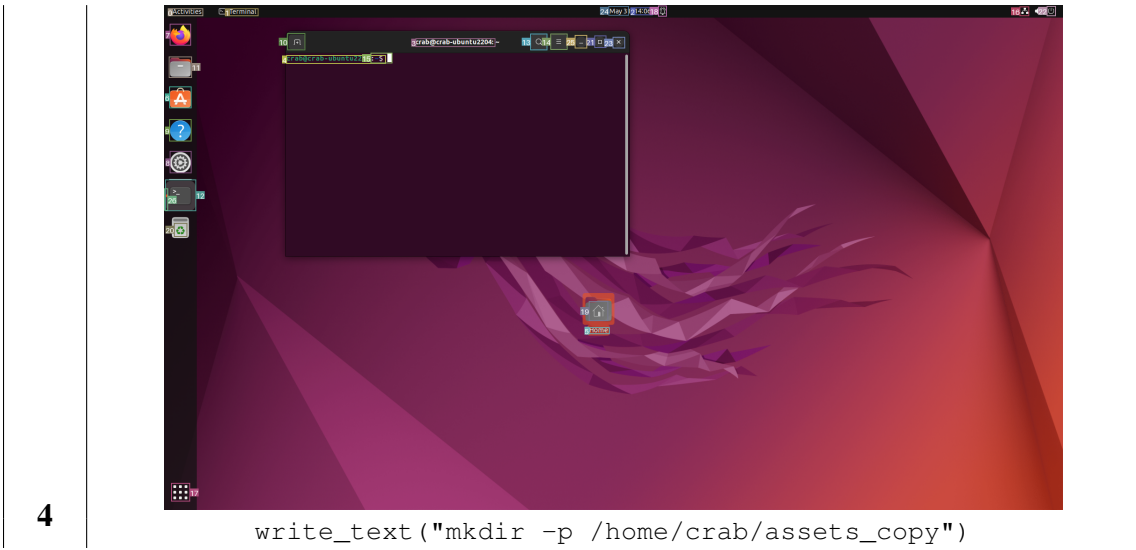
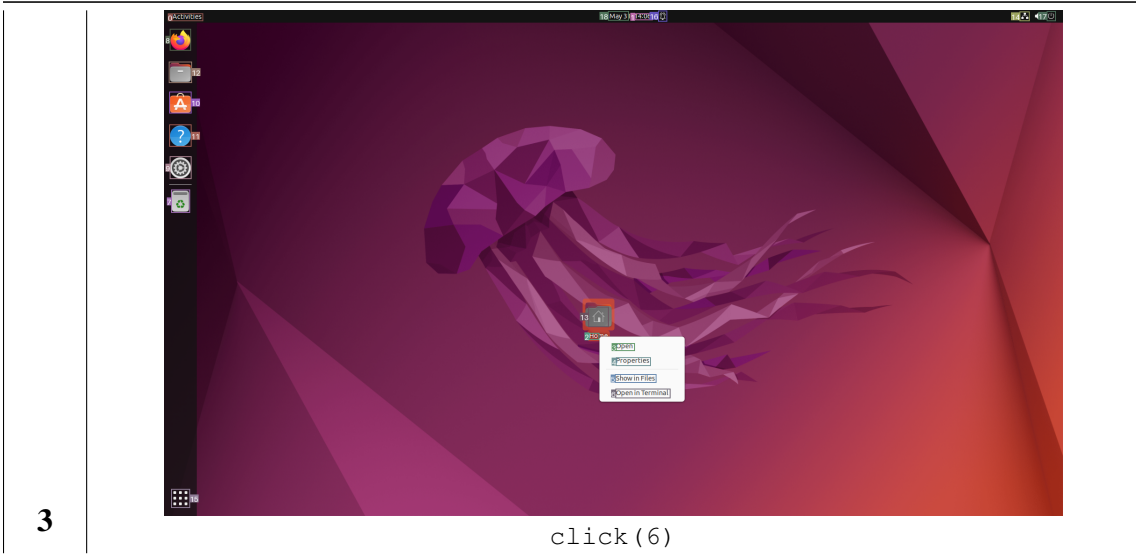
Table 9: **Ubuntu task case with GPT-4o (Single):** Create a new directory "/home/crab/assets_copy" and copy all files with the specified "txt" extension from "/home/crab/assets" to the directory "/home/crab/assets_copy".

Step	Agent Observation and Action
0	 <p data-bbox="829 926 951 953">click (2)</p>
1-2	 <p data-bbox="781 1480 1000 1535">click (2) right_click (2)</p>

Continued on next page

Table 9 continued from previous page

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997



Continued on next page

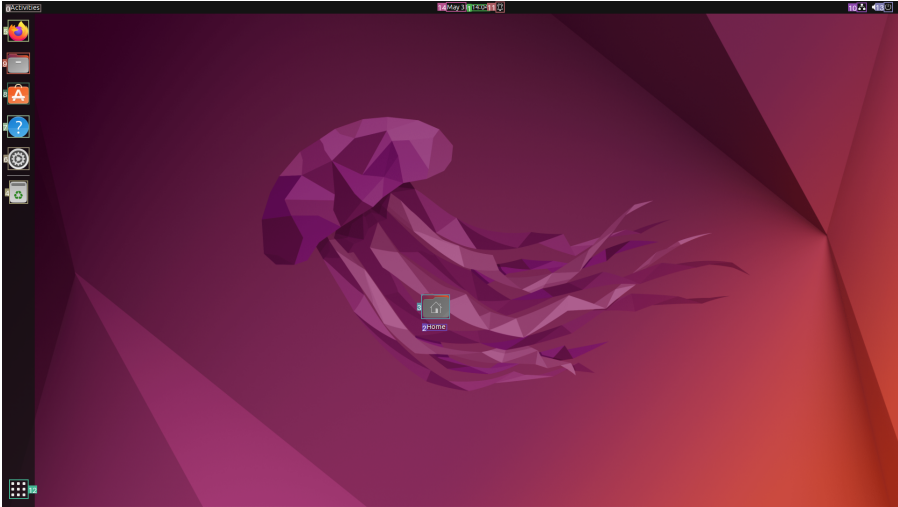
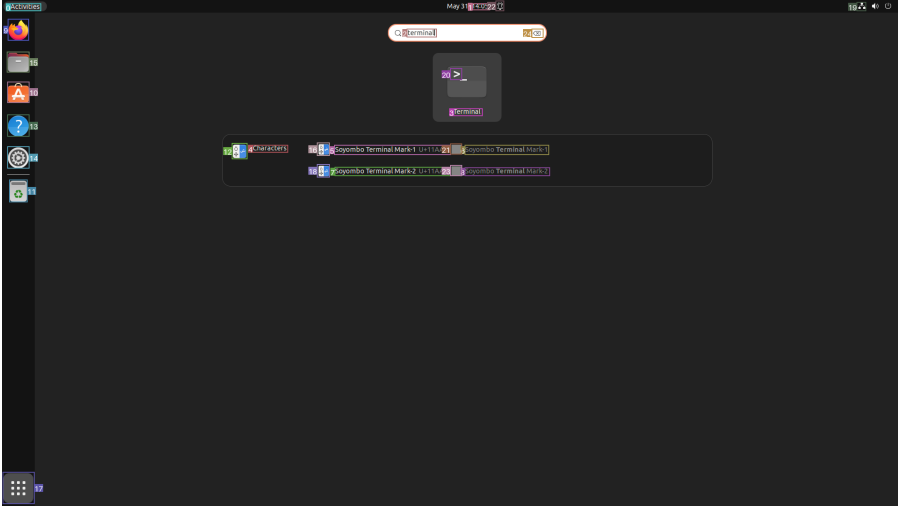
Table 9 continued from previous page

1998		
1999		key_press("enter")
2000	5-6	write_text("cp /home/crab/assets/*.txt /home/crab/assets_copy/")
2001		
2002		
2003		
2004		
2005		
2006		
2007		
2008		
2009		
2010		
2011		
2012		
2013		
2014		
2015		
2016		
2017	7	key_press("enter")

Completed Nodes / Total Nodes: 2 / 2 Termination Reason: Success

2051

2052 Table 10: **Ubuntu task case with GPT-4 Turbo (Single):** Create a new
 2053 directory `"/home/crab/assets_copy"` and copy all files with the specified
 2054 `"txt"` extension from `"/home/crab/assets"` to the directory `"/home/crab/
 2055 assets_copy"`.
 2056

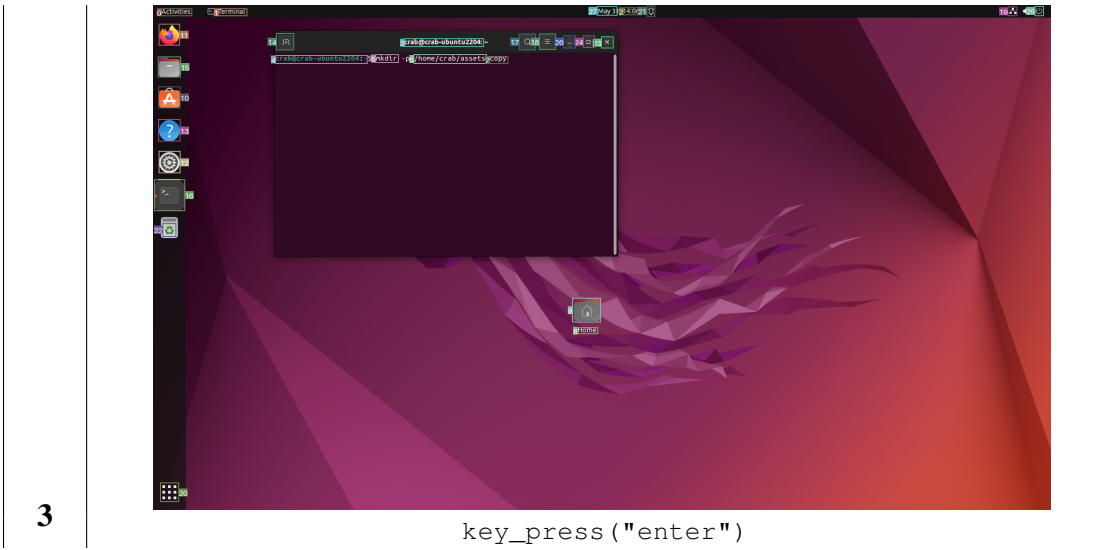
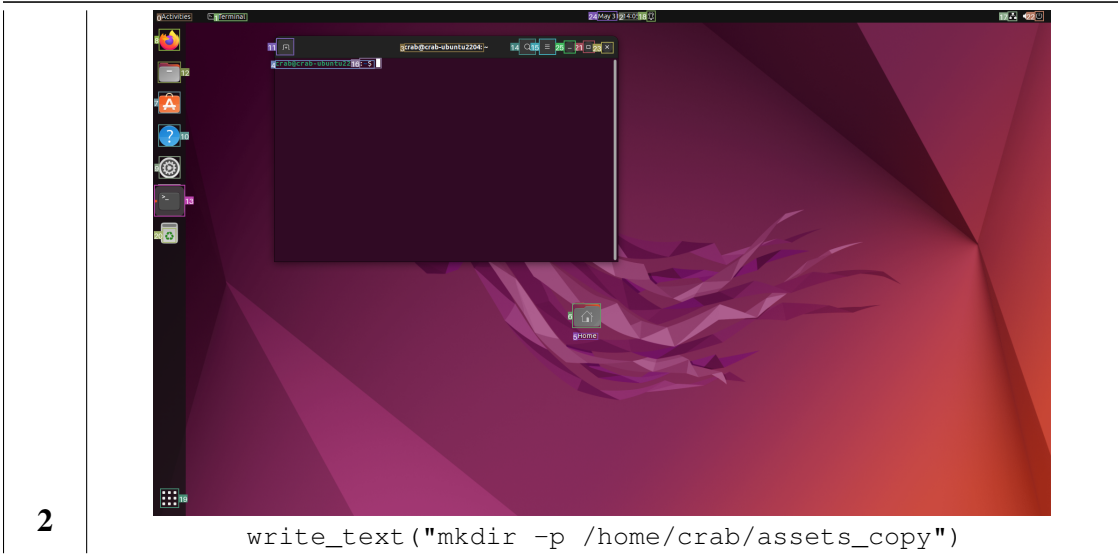
Step	Agent Observation and Action
2057 2058 2059 2060 2061 2062 2063 2064 2065 2066 2067 2068 2069 2070 2071 2072 2073 2074 2075	 <p data-bbox="646 926 1117 957">search_application("terminal")</p>
2076 2077 2078 2079 2080 2081 2082 2083 2084 2085 2086 2087 2088 2089 2090 2091 2092	 <p data-bbox="808 1480 954 1512">click(20)</p>

Continued on next page

2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

Table 10 continued from previous page

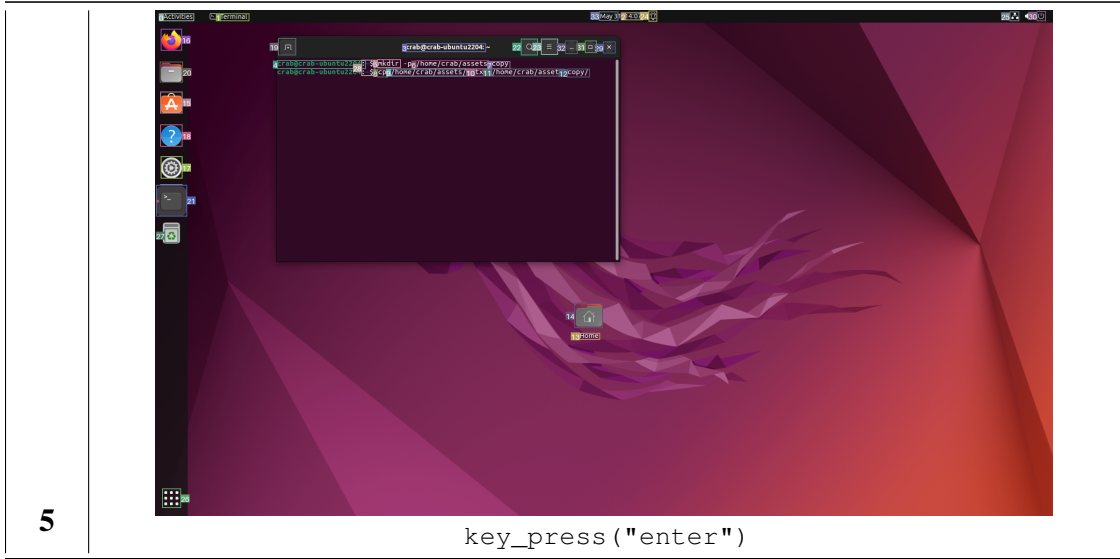
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159



Continued on next page

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

Table 10 continued from previous page



`key_press("enter")`

Completed Nodes / Total Nodes: 2 / 2 Termination Reason: Success