# 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 6 advanced MLMs using different single and multi-agent system configurations on this benchmark. The experimental results demonstrate that the single agent with GPT-40 achieves the best completion ratio of 38.01%.

025

004

010 011

012

013

014

015

016

017

018

019

021

### 1 INTRODUCTION

027 028 029

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 040 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 041 benchmarks that interact with the environments through pre-collected observation data from system 042 environments (Sun et al.; Mialon et al.; Deng et al., 2023) fail to capture the dynamic nature of real-043 world scenarios without interactive exploration where data and conditions can change unpredictably. 044 Second, existing benchmarks are typically evaluated on a single platform, either Web, Android, or 045 Desktop OS (Shi et al., 2017; Xing et al.; Xie et al.). However, the practical applications usually 046 involve tasks that span multiple platforms. For example, using a smartphone to take a photo and 047 sending it to a desktop for editing with a graphics editor is a common real-world task across multiple 048 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 051 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 052 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 environments, continuously updating and outputting the task completion metrics throughout the workflow.

079

073

074

075

076

077

We propose a benchmark that closely mirrors real-world situations and an evaluation method that more 081 accurately reflects an agent's performance on complex tasks. To this end, we introduce CRAB, a novel **CR**oss-environment Agent Benchmark framework. CRAB provides a comprehensive framework for 083 evaluating cross-environment tasks in interactive environments, where the agent needs to operate 084 simultaneously across various devices and platforms, adapting to varied system conditions to complete 085 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 087 method called graph evaluator. Unlike traditional goal-based and trajectory-based evaluation, our 088 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 089 each is considered a node in the graph evaluator. The graph structure describes the sequential and 090 parallel relationships between the sub-goals. Therefore, it offers fine-grained metrics similar to 091 trajectory-based evaluations while accommodating multiple valid pathways to a solution, making it 092 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 094 task construction method called **sub-task composition**. Combining multiple sub-tasks in a graph with 095 task targets allows for efficient construction of various cross-environment tasks with corresponding 096 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 098 CRAB with existing agent benchmark frameworks.

099 Based on CRAB framework, we propose a benchmark Crab Benchmark-v0 with two cooperated 100 environments that include an Android emulator and an Ubuntu desktop virtual machine. We have 101 developed a total of 120 real-world tasks. These tasks address a wide array of common real-world 102 applications and tools, including but not limited to calendars, email, maps, web browsers, and 103 terminals, and facilitate common interactions between smartphones and desktops. Considerable time 104 has been invested in verifying the accuracy and comprehensiveness of the instructions for subtasks, 105 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 106 collaborative scenarios. We test 6 popular MLMs, including GPT-4 Turbo, GPT-4o, Claude 3 Pro, 107 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 109 of each framework: Interactive Environment indicates the presence of either interactive environments 110 or static datasets; *Multimodal Observation* specifies the availability of vision-based observations (e.g. 111 screenshots); Cross-platform denotes support for multiple operating systems or platforms; Evaluation 112 describes the evaluation metrics, categorized as *Goal-based* (checking environment state according 113 solely on the final goal), Trajectory-based (comparing agent action trajectory with a gold actions 114 sequence), Multiple (varied across tasks), Intermediate-reward (combines multiple signals with three 115 strategies: Conjunctive Evaluation, Disjunctive Evaluation, and Order Constraint), LLM-as-a-Judge 116 (Zheng et al., 2023), or *Graph-based* (a DAG with each node as an intermediate checkpoint); *Task* 117 *Construction* shows the task construction method, including *Handmade* (handcrafted by human), 118 LLM-inspired (using LLM to generate task drafts but still verified and annotated by human), Template 119 (generated by filling in the blanks in task templates), or *Sub-task Composition* (composing multiple sub-tasks to construct tasks and evaluators). 120

	Interactive Environment	Multimodal Observation	Cross- platform	Evaluation	Task Construction	# of apps or websites
MINIWOB++ (Shi et al., 2017)	Web	1	×	Goal-based	Handmade	1
WEBSHOP (Yao et al., 2022)	Web	1	×	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	<ul> <li>Image: A set of the set of the</li></ul>	×	Goal-based	Template	6
OMNIACT (Kapoor et al.)	×	×	×	Trajectory-based	Handmade	60+
VWEBARENA (Koh et al.)	Web	1	×	Goal-based	Template	4
ANDROIDARENA (Xing et al.)	Android	1	×	Trajectory-based	LLM-inspired	9
OSWORLD (Xie et al.)	Linux / Windows	1	×	Goal-based	Template	9
MOBILE-ENV (Zhang et al., 2024a)	Android	<ul> <li>Image: A second s</li></ul>	×	Intermediate-reward	Template	13
GUI-WORLD (Chen et al., 2024a)	×	<ul> <li>Image: A second s</li></ul>	×	LLM-as-a-Judge	LLM-inspired	not present
ANDROIDWORLD (Rawles et al., 2024)	Android	<ul> <li>Image: A second s</li></ul>	×	Goal-based	Template	20
WAA (Bonatti et al., 2024)	Windows	<ul> <li>Image: A second s</li></ul>	×	Goal-based	Handmade	6
CRAB	Linux & Android	<ul> <li>Image: A second s</li></ul>	1	Graph-based	Sub-task Composition	25

133

134 135 136

137

138

139

and multi-agent systems, totaling 12 different agent settings in our benchmarks. The experimental results show that the single agent structure with GPT-40 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.

140 141 142

143

### 2 RELATED WORK

144 Leveraging LLMs as reasoning units has become an effective approach (Wang et al., 2024b; Huang 145 et al., 2022; Xi et al.) for building autonomous agents, including embodied agents (Wang et al., a; 146 Song et al., 2023; Chen et al., 2023), social simulations (Park et al., 2023; Lin et al., 2023), web 147 navigation (Lù et al.), game playing (Lan et al., 2023; Tan et al., 2024), office assistants Li et al. 148 (2024b), and code generation (Zhang et al., 2023). Specifically, some works apply LLMs to the 149 planning of embodied agents in complex environments (Wang et al., a; Song et al., 2023; Chen 150 et al., 2023). Others focus on simulating human behaviors and social communication by harnessing 151 LLMs' remarkable human-like understanding and generation capabilities (Park et al., 2023; Lin et al., 152 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., 153 2024; Wang et al., 2024a). In another approach, several studies have expanded the capacities of 154 agents by incorporating multimodal understanding, enabling agents to process diverse modalities of 155 input data such as images and text (Hong et al.; Liu et al., a; Furuta et al., 2024; Chen et al., 2024b). 156

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. 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

162 environment to simulate sufficiently complex web tasks. Several works (Koh et al.; He et al., 2024) 163 further broaden the web tasks, considering the visual tasks to build the benchmark for multi-modal 164 autonomous agents. SWEBench (Jimenez et al.) builds a benchmark based on the Github, focusing 165 on the coding capacity of understanding and solving issues. AgentBench (Liu et al., 2024) expands 166 the scope of agent applications within the domain of computer interaction tasks and encompasses the examination of these tasks across multiple complex environments. OMNIACT (Kapoor et al.) 167 incorporates the visual information of OS screen UI via segmentation and corresponding tagging, 168 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 170 meticulously delineating a standardized format for both the environment and the evaluation process. 171 WindowsAgentArena(Bonatti et al., 2024) focuses on the simulation of windows environment, 172 proposes a challenging set of windows-oriented task, gives a trustful evaluation for the popular 173 environment. 174

Current studies also focus on control tasks in mobile systems. MetaGUI (Sun et al.) divides the 175 mobile system control tasks into dialogues and GUI operation traces, collecting GUI traces based on 176 the collected dialogues. AITW (Rawles et al., 2023) produces a large dataset upon a large dataset 177 of real-world scenarios, and builds challenging multi-steps tasks based on the annotated single-step 178 tasks as a two-stage manner. MobileAgent (Wang et al., b) proposes tasks based on Ant Intelligent 179 Assistant(AIA) system, which integrates Standing Operating Procedure(SOP) information for the creation of subtasks. AITZ (Zhang et al., 2024b) constructs datasets with Chain-of-Thought (COT) 181 considerations, adding semantic annotations according to visual models at each step, and developing 182 the operational procedure for selected tasks. Mobile Agent Bench (Wang et al., 2024c) collects app 183 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 184 a fully functional environment for the Android system and provides a robust and reliable evaluation 185 of the agent's capacity in Android-oriented tasks. Mobile-ENV (Zhang et al., 2024a) introduces an intermediate reward mechanism where the environment generates signals based on its state. These 187 signals are combined into an intermediate reward using three types of aggregation operators. The 188 motivation behind this approach aligns closely with the problem addressed by our proposed graph 189 evaluator, but it relies on a tree structure with multiple relationships that can increase complexity 190 for annotators, potentially limiting dataset scalability. GUI-World (Chen et al., 2024a) contributes 191 a large video dataset for GUI automation and trains a new VideoLLM for UI tasks. However, the 192 evaluation method used by GUI-World, LLM-as-a-Judger (Zheng et al., 2023), may lack the precision 193 and consistency offered by rule-based evaluation systems.

194 195 196

197

199

### 3 DEFINITIONS

### 3.1 PROBLEM FORMULATION

200 Consider autonomous agents performing a task on a digital device (i.e. desktop computer). Such a 201 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 202 environment. Formally, this environment is defined as a reward-free Partially Observable Markov 203 Decision Process (POMDP), denoted by the tuple M := (S, A, T, O), where S represents the state 204 space,  $\mathcal{A}$  the action space,  $\mathcal{T}: \mathcal{S} \times \mathcal{A} \to \mathcal{S}$  the transition function, and  $\mathcal{O}$  the observation space. 205 Considering the collaborative nature of multiple devices in real-world scenarios, we can combine 206 multiple environments into a set  $M = M_1, M_2, ..., M_n$ , where n is the number of environments and 207 each environment  $M_j = (S_j, A_j, T_j, O_j)$ . We define a task that requires operations across multiple 208 environments as a **cross-environment task**. This task is formalized as a tuple  $(\mathbf{M}, I, R)$ , in which 209 M is the environment set, I is the task objective in the form of natural language instructions, and 210 R is the reward function of the task. An **agent system**, designed to complete a task represented 211 by an instruction I, can be modeled as a policy  $\pi((m,a) \mid (I,H,o_1,\ldots,o_n))$ , which defines 212 the probability of taking action a in environment m when receiving observation  $(o_1, ..., o_n)$  from 213 environment  $(M_1, ..., M_n)$  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 214 system is composed of either a single agent responsible for all planning, reasoning, and action-taking 215 or multiple agents connected through a communication workflow to collaborate.

### 216 3.2 GRAPH OF TASK DECOMPOSITION

218 Decomposing a complex task into several sim-219 pler sub-tasks has been proved to be an effec-220 tive prompting method for LLMs (Khot et al., 2023). Some studies represent sub-tasks in a graph structure. For instance, PLaG (Lin et al.) 222 uses a graph-based structure to enhance plan rea-223 soning within LLMs, while DyVal (Zhu et al., 224 2024) employs directed acyclic graphs (DAGs) 225 to facilitate dynamic evaluation of LLMs. By 226 introducing this concept into the realm of bench-227 marks, naturally, decomposing a complex task 228 into sub-tasks that have both sequential and par-229 allel connections forms a DAG. Therefore, we 230 introduce the Graph of Decomposed Tasks 231 (GDT), which provides a new task decomposition method representing decomposed sub-tasks 232 within a DAG structure. In GDT, each node is a 233 sub-task, formalized as a tuple (m, i, r), where 234



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

m specifies the environment in which the sub-task is performed, *i* provides the natural language instruction, and *r* represents the reward function. This function evaluates the state of *m* and outputs a boolean value to determine if the sub-task is completed. The edges within GDT represent the sequential relationship between sub-tasks. An example GDT is shown in Fig. 2.

4 THE CRAB FRAMEWORK

### . . .

### 4.1 CROSS-ENVIRONMENT AGENT INTERACTION

244 Compared to single-environment tasks, cross-environment tasks offer two main advantages for 245 benchmarking agents. First, cross-environment tasks reflect real-world scenarios where humans 246 use multiple devices simultaneously to accomplish tasks. Second, these tasks require sophisticated 247 message processing and information transfer between environments. Such tasks demand that the agent 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. Through this approach, CRAB can adapt to any platform or modality, from devices to 252 applications like browsers, by defining a few interactive functions. Implementation details are in the 253 Appendix A.3. 254

255 256

257

239 240

241 242

243

### 4.2 GRAPH EVALUATOR

To assess the capabilities of MLM agents, most benchmarks (Shi et al., 2017; Deng et al., 2023; Koh 258 et al.; Zhou et al.) evaluate agents based on solely the final states of the environment after agent 259 operations. Typically, they only judge whether the final goal is success or fail. However, this approach 260 does not capture incremental progress made by the agents. For instance, consider two agents tasked 261 with installing a new application on a computer: agent a successfully downloads the installer but fails 262 during the installation process, whereas agent b does not even try to find the installer. Despite Agent 263 a making more progress, both are deemed failures under the goal-based evaluation system, resulting 264 in an unfair assessment of their performance. An alternative method, Trajectory-based Matching 265 (Xing et al.; Kapoor et al.), abandons state-based evaluation and instead compares the agent's actions 266 against a predefined gold action sequence for each task, giving nuanced metrics. Nevertheless, this 267 method faces challenges in real-world systems where tasks may have multiple valid execution paths. For example, copying a file can be accomplished using either a file manager or the command line. 268 Inspired by the "decomposing" idea from GDT (Sec. 3.2), we propose a novel integrated approach, 269 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 node completion triggers successor nodes to activate and verify the state. This cycle repeats until no new nodes activate, showing that the system's task sequence aligns with the current state of the 276 environment. Unlike trajectory-based methods, the Graph Evaluator focuses on key states rather than specific actions, allowing agents flexibility in execution. For instance, in a file-editing task, the 278 evaluator checks if the file is edited, regardless of whether a CLI or GUI editor is used. This ensures 279 mandatory steps are completed while accommodating diverse execution paths. 280

### 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 three metrics aiming at assessing both performance and efficiency of agents, leveraging the detailed 286 287 sub-task status provided by the graph evaluator. Specifically, the **Completion Ratio** (CR) measures the proportion of completed sub-task nodes relative to the total nodes in the graph, calculated as 288  $C \mid N$ , where C is the number of completed nodes and N is the total number of nodes. This 289 metric offers a straightforward measure of an agent's progress on a given task. The Execution 290 Efficiency (EE), calculated as CR / A, where A denotes the count of executed actions. It evaluates 291 how efficiently actions are executed relative to the completion of nodes, reflecting the agent's task 292 execution efficiency. Lastly, the Cost Efficiency (CE), calculated as CR / T, where T is the total 293 number of model tokens used, evaluates the efficiency of resource consuming by the agent.

294 295

296

281

282

### 4.4 TASK AND EVALUATOR CONSTRUCTION

297 Despite the graph evaluator offering detailed evaluations, its creation is complex, requiring task 298 decomposition into sub-tasks with well-defined graph structures and expert involvement. To ease task 299 and evaluator creation, we propose building GDTs by sub-tasks, addressing two main challenges: (1) the need for manual creation of sub-tasks and (2) the complexity of modeling sequential and parallel 300 relationships between them. A template-based approach is commonly used to address the first issue 301 by generating a large number of tasks efficiently. To tackle the second challenge, we employ the 302 message transferring concept (Sec. 4.1). Specifically, if a sub-task  $\alpha$  produces an output message that 303 serves as an input for another sub-task  $\beta$ , then  $\alpha$  can be considered a legitimate prerequisite of  $\beta$ , 304 allowing us to connect  $\alpha$  and  $\beta$  with an directed edge in the GDT. To further refine our approach, 305 we introduce a *sub-task template* structure. Each sub-task is described using a natural language 306 instruction template that includes several replaceable input attributes. The types of each input attribute 307 and the task output should be defined carefully. To generate a GDT, input attributes can be filled 308 with either a hand-crafted value corresponding to their type or linked to a task with the same output 309 type as the input type. From the evaluator's perspective, each sub-task template is linked to an evaluator generator that uses the input attribute value to generate evaluator subgraphs. Once a GDT is 310 constructed, the graph evaluator is created by interlinking each subgraph. We follow the principle 311 that each subtask should **do one thing within a single environment**, with clearly defined inputs and 312 outputs that enable seamless integration with other tasks. For example, downloading a file from a 313 URL to a file path is a well-defined subtask: it accepts a URL as input and outputs the file's contents. 314

Task descriptions are initially generated by GPT-4 from sub-task prompts and refined by human
 reviewers. This approach, unlike naive templates, allows for a more detailed and scalable task compo sition. Our method automates graph evaluator generation, relieving users of coding requirements and
 making the system accessible to a broader audience.

319

321

### 320 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

324 environments in a reproducible and standalone manner and utilize snapshots to ensure a consistent 325 initial state for all environments. The observation space consists solely of the current system screen 326 for both environments, captured in image format at each step of the agent's interaction. We employ 327 the Set-of-Marks visual prompt method (Yang et al., a) to label each interactive element on the screen. 328 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 329 through EasyOCR<sup>1</sup> to detect and label interactive text elements. Each detected item is assigned a 330 unique integer ID, facilitating reference within the action space. The action spaces for Ubuntu and 331 332 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 333 action to search for applications. For Android, the action set includes tapping actions, a text action, 334 a physical button action, and an action to open the app drawer. Additionally, we introduce three 335 environment-irrelevant actions: completing the task, submitting an answer and waiting. Detailed 336 descriptions for the environment implementation are shown in Appendix A.2. 337

Tasks. We meticulously construct 17 sub-task templates for the Android environment and 19 sub-task 338 templates for the Ubuntu environment. The Ubuntu templates encompass a variety of tasks such as 339 Command Line Interface (CLI) operations, file system management, search engine usage, desktop 340 configurations, and map navigation. Conversely, the Android sub-task templates are primarily focused 341 on the storage and transmission of messages via various applications. Each sub-task template is 342 linked to a graph evaluator consisting of one to four nodes. Each sub-task are its graph evaluator is 343 verified by at least two related field experts. We make sure that all tasks are reachable by human. We 344 generate 104 tasks by sub-task composition and make 16 tasks by hand to include more complex 345 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 346 more complex than those in other benchmarks, which naturally requires more time for design and 347 experimentation. A single sub-task in our benchmark might involves multiple operations across 348 several applications, unlike prior works where most tasks often focus on solving problems within 349 a single application. With multiple applications nature combined with the scalability of our task 350 composition and graph evaluator, our tasks are sufficiently challenging to test an agent's performance 351 across different applications and scenarios, thereby effectively assessing its generalization ability. The 352 format and the applications covered by the dataset are shown in Appendix A.4 and A.5, respectively. 353

**Evaluators.** To assess the intermediate states of sub-tasks as described in Sec. 4.2, we have imple-354 mented a comprehensive suite of execution-based evaluators. These evaluators retrieve and assess 355 specific current states, such as the edited content of a file or a modified setting, thereby determining 356 the successful completion of a sub-task. For each evaluator, input attributes are carefully selected to 357 interpret software information or system settings relevant to the scenario defined for the sub-task. 358 For instance, evaluators use file paths before and after edits as input parameters to verify the com-359 pletion of file editing sub-tasks. Specifically, for sub-tasks on the Android platform, we incorporate 360 XML-based evaluators (Xing et al.). We dump UI layout as XML path and verify whether the UI 361 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 362 acquiring necessary state information through conventional APIs is challenging. Image matching 363 offers fine-grained visual correspondences by comparing keypoint features between images, allowing 364 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 366 specified images, text content, and layouts-tasks for which trivial evaluation methods are lacking. 367 We utilize EasyOCR<sup>1</sup> and XFeat<sup>2</sup> as our primary tools for OCR and image matching. For tasks with 368 real-time characteristics that may change over time, we implement crawler scripts to capture dynamic 369 values at the moment of evaluation. These values are then compared with the results achieved by the 370 agent upon task completion. We have a total of 59 evaluator functions with different types. Each task 371 has 4.2 evaluators in average of the whole dataset. 372

- 373
- 374
- 375

<sup>376</sup> 377

<sup>&</sup>lt;sup>1</sup>https://github.com/JaidedAI/EasyOCR

<sup>&</sup>lt;sup>2</sup>https://github.com/verlab/accelerated\_features

### <sup>378</sup> 6 EXPERIMENTS

### 3806.1BASELINE AGENT SYSTEM381

382 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 384 input, as the system provides both screenshots and text instructions as prompts; (2) Handle multi-turn 385 conversations, as most tasks require the agent to take multiple actions, necessitating the storage of 386 history messages in its context; (3) Generate structured output through function calling, ensuring the 387 proper use of provided actions with type-correct parameters. However, most open source models do 388 not provide explicit function calling feature, we let these models generate structured JSON output to 389 simulate the function calling behavior. 390

We selected 4 commercial and 2 open source MLMs that meet these criteria for our experiments: 391 GPT-4o (gpt-4o-2024-05-13) (OpenAI, 2024), GPT-4 Turbo (gpt-4-turbo-2024-04-09) (Achiam 392 et al.), Gemini 1.5 Pro (May 2024 version) (Reid et al.), Claude 3 Opus (claude-3-opus-20240229) 393 (Anthropic, Year), Pixtral-12B (Pixtral-12B-2409)<sup>3</sup>, and LLaVA-OneVision-72B (Ilava-onevision-394 qwen2-72b-ov-chat) (Li et al., 2024a). These models serve as the backend models for our agents. 395 Specifically, We use function calling feature in the four commercial models and JSON output in the 396 two open source models that do not support function calling. Since the JSON output setting uses 397 different prompts from the other, we employ a GPT-40 agent without function calling as the control 398 group to the open source models. 399

Beyond the MLM backend, the structure of agent systems also influences overall performance. To 400 examine how different multi-agent structures impact performance, we design three agent system 401 structures: single agent, multi-agent by functionality, and multi-agent by environment. In 402 the **single agent** structure, one agent manages all responsibilities, including observation analysis, 403 planning, reasoning, and format the output action. The **multi-agent by functionality** structure splits 404 tasks between a main agent, responsible for analysis and planning, and a tool agent that translates 405 instructions into actions without accessing environmental observations. This division allows the 406 main agent to concentrate on high-level tasks without managing functional call formats. Meanwhile, 407 in the multi-agent by environment setup, responsibilities are further distributed. A main agent 408 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 409 their respective environments. 410

411 For all models, we utilized the default API parameters and retained two turns of historical messages 412 to ensure messages do not exceed the context window. The interaction turns are limited to 15 and 413 the task will be terminated when reaching max turns. The agent can also terminate the task ahead 414 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 415 experiment, we deployed four cloud machines cloned from the same disk image to ensure a consistent 416 environment for all agents. Running a single agent setting in the benchmark requires at least 30 hours 417 to complete on one machine. Evaluation duration depends on the agent system, API response time, 418 and task steps. Single-agent systems average 10 to 20 seconds per step, while multi-agent systems 419 take 20 to 40 seconds. 420

- 421
- 422

431

### 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.

<sup>&</sup>lt;sup>3</sup>https://mistral.ai/news/pixtral-12b/

Table 2: Evaluation results on Crab Benchmark-v0. The Model column identifies the backend 433 masked language models (MLMs) used. The Structure column describes the configuration of the 434 agent system: Single means single agent; By Func is multi-agent by functionality; By Env indicates 435 *multi-agent by environment.* We provide traditional metric of *Success Rate* (SR) alongside newly 436 introduced metrics: Completion Ratio (CR), Execution Efficiency (EE), and Cost Efficiency (CE). 437 Note that Gemini 1.5 Pro has an invalid CE because the Gemini API does not support retrieving 438 token counts at the start time of experiments. The *Termination Reason* shows the ratio of reasons 439 why the agent is terminated when the task is not success. False Completion (FC) indicates that the 440 agent believes it has completed the task, but it actually has not; Reach Step Limit (RSL) means the 441 agent has reached the step limit but has not completed the task; *Invalid Action* (IA) refers to the agent 442 producing outputs that do not follow instructions, which may include invalid formats, nonexistent 443 actions, or invalid action parameters.

4	Agent syst	tem		Me	etrics		Tern	nination Rea	ason
	Model	Structure	<b>SR(%)</b> ↑	$  CR(\%) \uparrow$	<b>EE(%)</b> ↑	<b>CE(%)</b> ↑	FC(%)	RSL(%)	IA(%)
	GPT-40	Single	14.17	<b>38.01</b>	<b>4.15</b>	$5.29 \times 10^{-4}$	8.33	55.83	21.67
	GPT-40	By Func	<b>15.00</b>	34.00	3.93	$5.31 \times 10^{-4}$	10.83	54.17	20.00
	GPT-4 TURBO GPT-4 TURBO	Single By Func	9.17	33.35 33.48	3.80 4.07	$   \begin{array}{r}     2.74 \times 10 \\     4.52 \times 10^{-4} \\     4.38 \times 10^{-4}   \end{array} $	8.33 10.83	48.33 65.00 40.00	17.50 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	$1.85 \times 10^{-4}$	10.00	57.50	29.17
	CLAUDE 3 OPUS	By Func	3.33	16.48	1.72	$1.77 \times 10^{-4}$	28.33	34.17	34.17
	GPT-40 w/o FC	Single	9.17	23.05	2.34	$3.93 \times 10^{-4}$	5.00	42.50	43.33
	Pixtral-12B	Single	0.83	9.50	0.75	$0.87 \times 10^{-4}$	0.83	75.83	22.50
	LLaVA-OV-72B	Single	0.83	6.64	0.52	$1.02 \times 10^{-4}$	12.50	71.67	15.00

459

460

461 efficiency, the GPT-4 series also demonstrates strong performance, with GPT-40 having a higher CE 462 value compared to GPT-4 Turbo, highlighting its cost-effectiveness. GPT-4o's performance drops 463 after disabling tool calling feature, primarily due to its higher Invalid Action rate, showing the effectiveness of tool calling in generating structured output. In open source models, Pixtral-12B, with 464 far fewer parameters, achieves a better CR compared to LLaVA-ov-72B, showcasing its efficiency. 465 Although the open-source models generally understand screenshots and generate step-by-step plans 466 correctly, they often fail to execute the correct actions according to the plan. Moreover, they do not 467 effectively analyze task completion through observation. Once an incorrect action is performed, they 468 tend to assume current step is success and proceed to the next step. 469

**Comparison of agent structures.** The performance of multi-agent structures on all backend MLMs 470 is slightly lower than that of single-agent structures, which is somewhat unconventional. Based on 471 the communication log, we find that multi-agent structures tend to experience information loss during 472 inter-agent communication, leading to misunderstandings among downstream agents. This increases 473 the likelihood of multi-agent structures taking invalid actions and incorrectly completing tasks. These 474 experiments demonstrate that the design of the communication protocol and selecting the appropriate 475 scenario are crucial for multi-agent systems. A detailed analysis is included in Appendix C.2. In 476 terms of efficiency, multi-agent structures require more chat rounds, which can consume more tokens, 477 resulting in a lower CE compared to single-agent settings.

478 **Comparison of platforms.** We have three types of tasks: Ubuntu, Android, and cross-environment. 479 The metrics for each type of task can reveal the model or structure preferences. As shown in Fig. 3, 480 the GPT-40 model demonstrates significantly better performance on Android and cross-platform 481 tasks compared to GPT-4 Turbo, which may indicate OpenAI's increased focus on mobile devices. 482 Additionally, models like Gemini, Claude, Pixtral, and LLaVa-OV perform better on Android devices 483 compared to the Ubuntu, likely due to less training on Linux desktop data, which makes it difficult for them to recognize desktop icons. While it does not fully represent agent performance, it's notable 484 that the two open-source models exhibit a low invalid action ratio but still fail to complete tasks. We 485 include further platform specific results in Appendix C.1.



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

Comparison of metrics. The completion ratio metric reveals a notable performance difference between models. For instance, even though GPT-40 with single agent structure and with mutli-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.

Key issues in solving cross-environment task. The benchmark pipeline's complexity makes it 514 difficult to identify universal issues across tasks and models. However, the challenges in cross-515 platform tasks are similar to those in single-platform settings. Key issues include **action space** 516 **discrepancies**, where diverse action spaces in cross-platform environments confuse single-agent 517 architectures but can be mitigated by multi-agent setups tailored to each platform; **limited context** 518 length, which prevents the ability to process entire history observations and becomes more severe for 519 cross-platform scenarios with increasing screenshots; coordinate grounding issues, where advanced 520 tools like GroundingDINO and OCR occasionally fail to detect all screen elements in too complicated 521 GUI observation; and **icon recognition failures**, where the backend model correctly plans the next 522 step but cannot accurately identify and interact with corresponding icons, even though the visual prompt detect them correctly. 523

524 525

526 527

505 506 507

508

509

510

511

512

513

### 7 CONCLUSION

We propose the CRAB framework, which introduces the cross-environment automatic task-performing 528 problem, featuring advanced graph-based task generation and evaluation methods that reduce manual 529 effort in task design and provide more dynamic and accurate agent assessments. Based on this 530 framework, we present Crab Benchmark-v0, a set of high-quality cross-environment tasks in smart-531 phone and desktop environments, equipped with advanced visual prompting techniques. We tested 532 various backend models and agent system structures on this dataset. The results reveal preferences for 533 different agent settings, demonstrating Crab Benchmark-v0's strong ability to distinguish MLMs and 534 autonomous agent systems. Despite our contribution to advancing cross-environment agent research, 535 there are still some limitations. The sub-tasks are built upon the original apps in the Ubuntu and 536 Android systems on Pixel devices, which limits the coverage of a wider range of applications. The 537 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, 538 prompts, and multi-agent structures, as well as improving the use of visual prompting methods within the benchmark.

### 540 REFERENCES

565

566

567

571

572

573 574

575

576

581

582

583

584

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
URL http://arxiv.org/abs/2303.08774.

- Anthropic. The claude 3 model family: Opus, sonnet, haiku. https://www-cdn.anthropic.
   com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model\_Card\_Claude\_3.
   pdf, Year.
- Fabrice Bellard. Qemu, a fast and portable dynamic translator. In USENIX annual technical conference, FREENIX Track, volume 41, pp. 10–5555. California, USA, 2005.
- Rogerio Bonatti, Dan Zhao, Francesco Bonacci, Dillon Dupont, Sara Abdali, Yinheng Li, Yadong Lu, Justin Wagle, Kazuhito Koishida, Arthur Bucker, Lawrence Jang, and Zack Hui. Windows agent arena: Evaluating multi-modal os agents at scale, 2024. URL https://arxiv.org/abs/2409.08264.
- Dongping Chen, Yue Huang, Siyuan Wu, Jingyu Tang, Liuyi Chen, Yilin Bai, Zhigang He, Chenlong
  Wang, Huichi Zhou, Yiqiang Li, Tianshuo Zhou, Yue Yu, Chujie Gao, Qihui Zhang, Yi Gui, Zhen
  Li, Yao Wan, Pan Zhou, Jianfeng Gao, and Lichao Sun. GUI-WORLD: A Dataset for GUI-oriented
  Multimodal LLM-based Agents, June 2024a.
- Wentong Chen, Junbo Cui, Jinyi Hu, Yujia Qin, Junjie Fang, Yue Zhao, Chongyi Wang, Jun Liu, Guirong Chen, Yupeng Huo, Yuan Yao, Yankai Lin, Zhiyuan Liu, and Maosong Sun. Guicourse:
   From general vision language models to versatile gui agents, 2024b. URL https://arxiv.org/abs/2406.11317.
  - Yongchao Chen, Jacob Arkin, Yang Zhang, Nicholas Roy, and Chuchu Fan. Autotamp: Autoregressive task and motion planning with llms as translators and checkers. *arXiv preprint arXiv:2306.06531*, 2023.
- Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong
   Wu. Seeclick: Harnessing gui grounding for advanced visual gui agents. URL http://arxiv. org/abs/2401.10935.
  - Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web, 2023.
  - Johan Edstedt, Qiyu Sun, Georg Bökman, Mårten Wadenbäck, and Michael Felsberg. RoMa: Robust Dense Feature Matching. 2024.
- Hiroki Furuta, Kuang-Huei Lee, Ofir Nachum, Yutaka Matsuo, Aleksandra Faust, Shixiang Shane
  Gu, and Izzeddin Gur. Multimodal web navigation with instruction-finetuned foundation models.
  In *The Twelfth International Conference on Learning Representations*, 2024. URL https:
  //openreview.net/forum?id=efFmBWioSc.
  - Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. Exploring network structure, dynamics, and function using networkx. In Gaël Varoquaux, Travis Vaught, and Jarrod Millman (eds.), *Proceedings of the 7th Python in Science Conference*, pp. 11–15.
- Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan, and Dong Yu. Webvoyager: Building an end-to-end web agent with large multimodal models, 2024.
- Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang,
   Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. Metagpt: Meta programming for multi-agent
   collaborative framework. *arXiv preprint arXiv:2308.00352*, 2023.
- Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan
   Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent: A visual language model for gui agents. URL https://arxiv.org/abs/2312.08914v2.

611

619

639

- 594 Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot 595 planners: Extracting actionable knowledge for embodied agents. In International Conference on 596 Machine Learning, pp. 9118–9147. PMLR, 2022. 597
- Hanwen Jiang, Arjun Karpur, Bingyi Cao, Qixing Huang, and Andre Araujo. Omniglue: Gener-598 alizable feature matching with foundation model guidance. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2024. 600
- 601 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik 602 Narasimhan. Swe-bench: Can language models resolve real-world github issues? URL http: 603 //arxiv.org/abs/2310.06770.
- Dongming Jin, Zhi Jin, Xiaohong Chen, and Chunhui Wang. Mare: Multi-agents collaboration 605 framework for requirements engineering, 2024. URL https://arxiv.org/abs/2405. 606 03256. 607
- 608 Raghav Kapoor, Yash Parag Butala, Melisa Russak, Jing Yu Koh, Kiran Kamble, Waseem Alshikh, 609 and Ruslan Salakhutdinov. Omniact: A dataset and benchmark for enabling multimodal generalist 610 autonomous agents for desktop and web. URL http://arxiv.org/abs/2402.17553.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and 612 Ashish Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. 613 In The Eleventh International Conference on Learning Representations, 2023. URL https: 614 //openreview.net/forum?id=\_nGgzQjzaRy. 615
- 616 Avi Kivity, Yaniv Kamay, Dor Laor, Uri Lublin, and Anthony Liguori. kvm: the linux virtual machine 617 monitor. In Proceedings of the Linux symposium, volume 1, pp. 225-230. Dttawa, Dntorio, Canada, 618 2007.
- Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham 620 Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating 621 multimodal agents on realistic visual web tasks. URL http://arxiv.org/abs/2401. 622 13649. 623
- 624 Yihuai Lan, Zhiqiang Hu, Lei Wang, Yang Wang, Deheng Ye, Peilin Zhao, Ee-Peng Lim, Hui Xiong, 625 and Hao Wang. Llm-based agent society investigation: Collaboration and confrontation in avalon gameplay. arXiv preprint arXiv:2310.14985, 2023. 626
- 627 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, 628 Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer, September 2024a. 629
- 630 Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 631 Camel: Communicative agents for" mind" exploration of large scale language model society. 2023. 632
- Hongxin Li, Jingran Su, Yuntao Chen, Qing Li, and ZHAO-XIANG ZHANG. Sheetcopilot: Bringing 633 software productivity to the next level through large language models. Advances in Neural 634 Information Processing Systems, 36, 2024b. 635
- 636 Fangru Lin, Emanuele La Malfa, Valentin Hofmann, Elle Michelle Yang, Anthony Cohn, and Janet B. 637 Pierrehumbert. Graph-enhanced large language models in asynchronous plan reasoning. URL 638 http://arxiv.org/abs/2402.02805.
- Jiaju Lin, Haoran Zhao, Aochi Zhang, Yiting Wu, Huqiuyue Ping, and Qin Chen. Agentsims: An 640 open-source sandbox for large language model evaluation. arXiv preprint arXiv:2308.04026, 2023.
- 642 Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng Li, Tianhe Ren, Xueyan Zou, Jianwei Yang, 643 Hang Su, Jun Zhu, Lei Zhang, Jianfeng Gao, and Chunyuan Li. Llava-plus: Learning to use tools 644 for creating multimodal agents, a. URL https://arxiv.org/abs/2311.05437v1. 645
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, 646 Hang Su, Jun Zhu, and Lei Zhang. Grounding dino: Marrying dino with grounded pre-training for 647 open-set object detection, b. URL https://arxiv.org/abs/2303.05499v4.

648 649 650 651 652	Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. Agentbench: Evaluating LLMs as agents. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=zAdUB0aCTQ.
653 654 655	Xing Han Lù, Zdeněk Kasner, and Siva Reddy. Weblinx: Real-world website navigation with multi-turn dialogue. URL https://arxiv.org/abs/2402.05930v1.
656 657 658	Michael M. McKerns, Leif Strand, Tim Sullivan, Alta Fang, and Michael A. G. Aivazis. Building a framework for predictive science. URL http://arxiv.org/abs/1202.1056.
659 660	Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia: A benchmark for general ai assistants. URL http://arxiv.org/abs/2311.12983.
661 662 663 664 665 666	Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Belle- mare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. 518(7540):529–533. ISSN 1476-4687. URL https://www.nature.com/articles/ nature14236.
668	<b>OpenAI. Gpt-4 omni.</b> https://openai.com/index/hello-gpt-4o/, 2024.
669 670 671	Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In <i>Proceedings of the 36th</i> <i>Annual ACM Symposium on User Interface Software and Technology</i> , pp. 1–22, 2023.
673 674 675	Guilherme Potje, Felipe Cadar, Andre Araujo, Renato Martins, and Erickson R Nascimento. Xfeat: Accelerated features for lightweight image matching. In <i>IEEE/CVF Conference on Computer</i> <i>Vision and Pattern Recognition (CVPR)</i> , 2024.
676 677	Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. Android in the wild: A large-scale dataset for android device control, 2023.
678 679 680 681 682	Christopher Rawles, Sarah Clinckemaillie, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Marybeth Fair, Alice Li, William Bishop, Wei Li, Folawiyo Campbell-Ajala, Daniel Toyama, Robert Berry, Divya Tyamagundlu, Timothy Lillicrap, and Oriana Riva. Androidworld: A dynamic benchmarking environment for autonomous agents, 2024. URL https://arxiv.org/abs/2405.14573.
683 684 685 686	Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. URL http: //arxiv.org/abs/2403.05530.
687 688 689 690 691	Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Percy Liang. World of bits: An open-domain platform for web-based agents. In Doina Precup and Yee Whye Teh (eds.), <i>Proceedings of the 34th International Conference on Machine Learning</i> , volume 70 of <i>Proceedings</i> <i>of Machine Learning Research</i> , pp. 3135–3144. PMLR, 06–11 Aug 2017. URL https:// proceedings.mlr.press/v70/shi17a.html.
692 693 694 695	Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 2998–3009, 2023.
696 697 698	Liangtai Sun, Xingyu Chen, Lu Chen, Tianle Dai, Zichen Zhu, and Kai Yu. Meta-gui: Towards multi- modal conversational agents on mobile gui. URL http://arxiv.org/abs/2205.11029.
699 700 701	Weihao Tan, Ziluo Ding, Wentao Zhang, Boyu Li, Bohan Zhou, Junpeng Yue, Haochong Xia, Jiechuan Jiang, Longtao Zheng, Xinrun Xu, Yifei Bi, Pengjie Gu, Xinrun Wang, Börje F. Karlsson, Bo An, and Zongqing Lu. Towards general computer control: A multimodal agent for red dead redemption ii as a case study, 2024.

702 703 704 705 706 707 708 709 710	Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P. Agapiou, Max Jader- berg, Alexander S. Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David Budden, Yury Sulsky, James Molloy, Tom L. Paine, Caglar Gulcehre, Ziyu Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom Schaul, Timothy Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver. Grand- 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.
711 712 713	Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models, a. URL http://arxiv.org/abs/2305.16291.
714 715 716	Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances large language model capabilities, 2024a. URL https://arxiv.org/abs/2406.04692.
717 718 719	Junyang Wang, Haiyang Xu, Jiabo Ye, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang. Mobile-agent: Autonomous multi-modal mobile device agent with visual perception, b. URL https://arxiv.org/abs/2401.16158v2.
720 721 722	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. <i>Frontiers of Computer Science</i> , 18(6):1–26, 2024b.
723 724 725 726	Luyuan Wang, Yongyu Deng, Yiwei Zha, Guodong Mao, Qinmin Wang, Tianchen Min, Wei Chen, and Shoufa Chen. Mobileagentbench: An efficient and user-friendly benchmark for mobile llm agents, 2024c. URL https://arxiv.org/abs/2406.08184.
727 728 729	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. <i>arXiv preprint arXiv:2308.08155</i> , 2023.
730 731 732 733 734 735	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongx- iang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. The rise and potential of large language model based agents: A survey. URL https://arxiv.org/abs/2309.07864v3.
736 737 738 739	Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. Osworld: Benchmarking multimodal agents for open- ended tasks in real computer environments. URL http://arxiv.org/abs/2404.07972.
740 741 742 743	Mingzhe Xing, Rongkai Zhang, Hui Xue, Qi Chen, Fan Yang, and Zhen Xiao. Understanding the weakness of large language model agents within a complex android environment. URL http://arxiv.org/abs/2402.06596.
744 745 746	Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v, a. URL http://arxiv.org/abs/2310.11441.
747 748 749	John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. Intercode: Standardizing and benchmarking interactive coding with execution feedback, b. URL http://arxiv.org/abs/2306.14898.
750 751 752 753	Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. <i>Advances in Neural Information Processing Systems</i> , 35:20744–20757, 2022.
754 755	Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang, Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qingwei Lin, Saravan Rajmohan, Dongmei Zhang, and Qi Zhang. Ufo: A ui-focused agent for windows os interaction, a. URL http://arxiv.org/abs/2402.07939.

756	Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu.
757	Appagent: Multimodal agents as smartphone users, b. URL http://arxiv.org/abs/2312.
758	13771.
759	

- Danyang Zhang, Zhennan Shen, Rui Xie, Situo Zhang, Tianbao Xie, Zihan Zhao, Siyuan Chen, Lu Chen, Hongshen Xu, Ruisheng Cao, and Kai Yu. Mobile-Env: Building Qualified Evaluation Benchmarks for LLM-GUI Interaction, June 2024a.
- Jiwen Zhang, Jihao Wu, Yihua Teng, Minghui Liao, Nuo Xu, Xiao Xiao, Zhongyu Wei, and Duyu
   Tang. Android in the zoo: Chain-of-action-thought for gui agents, 2024b.
- 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.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
  Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
  Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena, December 2023.
- Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web environment for building autonomous agents. URL http://arxiv.org/abs/2307. 13854.
- Kaijie Zhu, Jiaao Chen, Jindong Wang, Neil Zhenqiang Gong, Diyi Yang, and Xing Xie. Dyval:
   Dynamic evaluation of large language models for reasoning tasks. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?
   id=gjfOL9z5Xr.

### 810 A BENCHMARK DETAIL

 Section A.2 introduces the implementation details and action space settings of the benchmark environments. Section A.3 describes the design logic and implementation of the CRAB framework. Section A.4 describes the our experiment settings in detail. Section A.5 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.<sup>4</sup> 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 DATASET STATISTICS

The applications in our task dataset along with the counts of tasks that utilize them is listed in
Table 3 and 4. The task dataset covers a wide range of applications across two platforms, primarily
focusing on daily life, programming, and office work scenarios. It is also worth noting that in our
task settings, a single task often involves two or more applications. On average, each task contains
1.84 applications, according to our statistics.

App Name	Description	# Tasks
Terminal	GNOME terminal emulator with command line tools (e.g., cat, wget).	40
Firefox	Web browser with various web Apps (e.g., Google Docs and Search).	35
File Manager	GNOME official file manager.	25
GIMP	GNU Image Manipulation Program, open-source raster graphics editor.	13
System Setting	GNOME system setting GUI application.	11
VSCode	Code editor.	8
LibreOffice Writer	Word processor.	8
LibreOffice Impress	Presentation program.	7
LibreOffice Calc	Spreadsheet program.	6
Vim	CLI text editor.	6
Slack	Team communication platform.	1

879

864

Table 4: Applications on	their tech counts in the	Android anvironment
Table 4. Additions and	I LHEIF LASK COUNTS III LHE	Апагоја епујгоншені.

880	App Name	Description	# Tasks
881	Google Map	Map application.	13
882	Google Calendar	Calendar application.	9
883	Gmail	Google mail service application.	7
884	Google Keep	Google note application.	6
885	Google Tasks	Google TO-DO list.	5
886	Messages	Android built-in message sending application.	5
887	Contacts	Android built-in contacts application.	5
888	Google Drive	Google Cloud Drive application.	4
880	Clock	Android built-in clock application.	2
005	Files	Android built-in file manager.	1
090	Settings	Android system setting.	1
891	Camera	Android built-in camera.	1
892	Google Docs	Google online word processor.	1
893	Phone	Android built-in phone calling application.	1

894 895

899

900 901

907

908

913

914

The distribution of node counts of graph evaluators per task is provided in Table 5. Our task dataset includes graphs ranging from 1 to 11 nodes. It is important to note that the number of nodes depends on the complexity of the task, with more complex tasks involving larger graphs.

A.2 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 Emulator<sup>5</sup>. Interaction with the
Ubuntu environment is facilitated using PyAutoGUI<sup>6</sup> and MSS<sup>7</sup>, 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)<sup>8</sup>. The detailed action space is described in Table 6.

A.3 FRAMEWORK DESIGN

909 CRAB offers a modular and extensible framework for evaluating agent performance in diverse tasks.
910 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

917 <sup>7</sup>https://github.com/BoboTiG/python-mss

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

<sup>915 &</sup>lt;sup>5</sup>https://developer.android.com/studio/run/emulator

<sup>916 &</sup>lt;sup>6</sup>https://github.com/asweigart/pyautogui

<sup>%</sup>https://developer.android.com/tools/adb

Table 5: Node count histogram. **# Nodes** # Tasks 

Table 6: 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.

0 | 3

27 29	Action Name (Parameters)	Description
)28 )29 )30 )31 )32 )32 )33 )34 )35	<pre>click(elem) right_click(elem) double_click(elem) write_text(text) press(key) hotkey(keys) scroll(direction) search_app(name)</pre>	Click on elem. Right-click on elem. Double-click on elem. Typing the specified text. Press a keyboard key. Press keyboard keys at the same time. Scrolls page up or down. Search for application with name in the system
)36 )37 )38 )39 )40 )41	<pre>tap(elem) long_tap(elem) swipe(elem, dire, dist) write_text(text) press(key) show_all_drawer()</pre>	Tap on elem. Press and hold elem. Swipe from elem in a specified direction and distance. Typing the specified text. Press a key, can be <i>home</i> or <i>back</i> . Show the app drawer to list installed applications.
942 943 944	<pre>submit(answer) complete() wait()</pre>	Submit answer if needed. State that a task is completed. Wait the environment to process

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.

- 969 A.4 CONFIGURATION BY MODULES
- 971 Building on the declarative and modular design of our framework, this section explains the configuration and potential extensibility of each module.

972 **Environment** The environments in CRAB are a combination of multiple different uses of actions 973 with some environment metadata, such as name and natural language description. In Crab Benchmark-974 v0, we use a computer desktop environment and a smartphone environment both based on virtual 975 machine technology. The computer desktop environment, named Ubuntu, is installed from an ISO 976 image of Ubuntu 22.04.4 LTS (Jammy Jellyfish) downloaded from the Ubuntu Official website<sup>3</sup>. 977 Necessary applications such as the LibreOffice suite (Writer, Calc, and Impress) and Slack are installed 978 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) 979 provided in Google Android Studio<sup>10</sup>. We install additional required applications such as *Keep Notes*, 980 Tasks, and Docs from Google Play. The descriptions of the two environments in Crab Benchmark-v0, 981 which are inserted in the agent prompts, are as follows: 982

- 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 997 provides an intuitive method to define actions through Python function. Here is an example of action 998 search application in the Ubuntu environment: 999

```
1000
      @action
      def search_application(name: str) -> None:
1001
           """Search an application name.
1002
1003
           For exmaple, if you want to open an application named "slack",
1004
           you can call search_application(name="slack"). You MUST use this
           action to search for applications.
           Args:
               name: the application name.
1008
           .....
1009
           pyautogui.hotkey("win", "a")
1010
           time.sleep(0.5)
           pyautogui.write(name)
1011
           time.sleep(0.5)
1012
1013
```

Listing 1: Define "search\_application" action.

We extract key information from the function through the *@action* decorator as following: 1015

- 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.

1021

1014

1016

1017

983

984

985

986

987

988

989

990

991

992

993

994

995

<sup>1023</sup> 1024 1025

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

<sup>&</sup>lt;sup>10</sup>https://developer.android.com/studio

1020	• <b>Parameters</b> : The parameters are the arguments that the functions accept, offering flexibility
1027	for the agent to control the environment. Typically, a set of parameters is defined, each
1028	consisting of a name, type, and a natural language description. Parameters are extracted from
1029	the function's parameters along with their type annotations. Additionally, parameter descrip-
1030	tions are extracted from the Args section in the docstring. In this example, there is only one
1031	parameter named name, with a type of str, and its description is the application
1032	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
1043
      def screenshot() -> str:
1044
           ""Capture the current screen as a screenshot."""
1045
          with mss() as sct:
1046
           # Capture raw pixels from the screen
1047
          sct_img = sct.grab(sct.monitors[1])
           # Convert to PNG format
1048
          png = tools.to_png(sct_img.rgb, sct_img.size)
1049
           # Encode to Base64 format for easier transmission
1050
          base64_img = base64.b64encode(png).decode("utf-8")
1051
          return base64_img
1052
```

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")
1063
      def check_text_in_current_window_name(text: str) -> bool:
1064
          try:
               out = subprocess.check_output(
1065
                   ["xdotool", "getwindowfocus", "getwindowname"], text=True
              ).strip()
1067
          except subprocess.CalledProcessError:
1068
               return False
1069
          return text in out
1070
```

1071

1033

1034

1035

1036

1038

1053

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.

1078

**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:

```
1080
        Task(
1081
             id="a3476778-e512-40ca-b1c0-d7aab0c7f18b",
1082
            description="Open \"Tasks\" app on Android, check the...",
1083
            evaluator=path_graph(
                 check_current_package_name("com.google.android.apps.tasks"),
1084
                 check_current_window_process("gnome-control-center"),
1085
                 check_color_scheme("prefer-dark"),
            ),
1088
                                          Listing 4: Define a task.
1089
1090
       In this model, each task is represented as an instance of the Task class, which is a subclass of
1091
       BaseModel in Pydantic<sup>11</sup> package. Each task is uniquely identified by an ID and described by a
1092
       detailed description. The evaluator component is structured as a graph evaluator, which integrates
1093
        multiple evaluative functions into a directed graph using the networkx<sup>12</sup> package. Each evaluator
1094
       within this graph must be appropriately parameterized to assess specific conditions relevant to the
1095
       task. For example, the task demonstrated aims to open the "Tasks" app on Android and perform
1096
       a series of verifications: it checks whether the correct Android app is opened, whether the current
1097
        focused window's process name is gnome-control-center, and whether the color scheme is
       set to dark.
1099
1100
       Sub-task The sub-task in CRAB is the unit component of in task construction. The following
1101
       example is a sub-task template that we used to easily generate sub-tasks:
1102
        SubTask(
1103
            id="0f589bf9-9b26-4581-8b78-2961b115ab49",
1104
            description="Open \"{file_path}\" using vim in a terminal, write \"{
            content}\", then save and exit vim.",
1105
            attribute_dict={"file_path": "file_path", "content": "message"},
1106
            output_type="file_path",
1107
            evaluator_generator=lambda file_path, content: path_graph(
1108
                 check_current_window_process("qnome-terminal-server"),
1109
                 is_process_open("vim"),
                 is_process_close("vim"),
1110
                 check_file_content(file_path, content),
1111
            ),
1112
        ),
1113
                                          Listing 5: Define a task.
1114
1115
       In this sub-task model, each sub-task is defined using a similar approach to the main task. The
1116
       attributes of the sub-task are outlined in an attribute_dict, which details the types and roles of
1117
       each attribute used in the sub-task's operations. The output_type field specifies the expected type
1118
       of output from the sub-task. The types reflected in attribute_dict and output_type, play
1119
       a critical role in determining the compatibility and sequential logic of compose multiple sub-tasks.
1120
       The evaluator for the sub-task is dynamically generated using a lambda function, which crafts an
1121
       evaluator sub-graph based on the sub-task's attributes.
1122
1123
        A.5 COMPOSED TASK FORMAT
1124
1125
       We use a JSON format to save the composed tasks, which includes the task ID, overall task description,
1126
        sub-tasks with their attribute values, and a graph structure represented in an adjacency list. The entire
1127
        task dataset is defined by the sub-task pool in Python code and the task composition JSON files
       categorized by task platform.
1128
1129
1130
             "description": "Combine Image 1 \"/home/crab/Pictures/cat.png\" and
1131
            Image 2 \"/home/crab/assets/campus.png\" using GIMP (GNU Image
1132
```

<sup>&</sup>lt;sup>11</sup>https://pydantic.dev/

<sup>&</sup>lt;sup>12</sup>https://networkx.org/

1134 Manipulation Program), placing Image 1 on the left side of Image 2, 1135 and save the combined image to \"/home/crab/Desktop/background.png\". 1136 Then, set this combined image as the screen background of the system .", 1137 "tasks": [ 1138 { 1139 "task": "4cf246ea-0a7f-43da-84b6-61d74a2699af", 1140 "attribute": { 1141 "image\_path\_1": "/home/crab/Pictures/cat.png", "image\_path\_2": "/home/crab/assets/campus.png" 1142 "output\_path": "/home/crab/Desktop/background.png" 1143 }, 1144 "output": "/home/crab/Desktop/background.png" 1145 }, 1146 "task": "a207ef38-b3b2-4c6c-a1e3-75c38162f5ba", 1147 "attribute": { 1148 "photo\_path": "/home/crab/Desktop/background.png" 1149 1150 "output": null 1151 } ], 1152 "adjlist": "0 1\n1", 1153 "id": "d3c917ff-406f-447a-87f5-b8d835cba750" 1154 1155 Listing 6: Define a composite task in JSON. 1156 1157 1158 В AGENT SYSTEM 1159 1160 **B.1** AGENT IMPLEMENTATION 1161 1162 In this section, we outline the implementation of the agents used in our experiments, which leverage 1163 advanced multimodal language models from OpenAI, Anthropic, and Google. Each agent is designed 1164 to function in multi-environment setups, interacting with various action spaces defined by different 1165 environments. 1166 1167 **General Framework** All agents share a common architecture but are tailored to the specific APIs 1168 and capabilities of each language model provider. 1169 1170 **Initialization** Each agent is initialized with several key parameters, including a description, an 1171 action space, the model type, maximum tokens, history message length, and an optional environment 1172 description. The initialization process involves: 1173 • Action Space Conversion: Actions defined for each environment are converted into a 1174 schema compatible with the respective API. This ensures that the actions can be correctly 1175 interpreted and executed by the language models. 1176 1177 • System Message Setup: Depending on whether the agent is configured for single or multiple 1178 environments, a system message is formatted to provide the model with context about the tasks and environments. 1179 1180 **Interaction (Chat Method)** The core functionality of each agent is encapsulated in its ability to 1181 interact with users through a chat method. This involves: 1182 1183 • **Content Parsing**: Input content is parsed and formatted to match the requirements of 1184 the respective API. This includes structuring user messages and any necessary contextual 1185 information. 1186 • Request Construction: The request payload is constructed, incorporating the system 1187

message, chat history, and the newly parsed user input.

- API Communication: The constructed request is sent to the appropriate API, which generates a response. The agents handle API-specific constraints such as rate limits and response formats.
  - **Response Handling**: The response from the API is processed to extract any tool calls suggested by the model. These are then appended to the chat history, maintaining a coherent conversation state.

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

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

**B.2** INTER-AGENT COMMUNICATION STRATEGIES

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

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

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

- 1222 B.3 AGENT PROMPT
- 1224 B.3.1 SINGLE AGENT

1225	
1226	Prompt
1227	You are a helpful assistant. Now you have to do a task as described below:
1228	**{task_description}**.
1229	You should never forget this task and always perform actions to achieve this task. And this is
1230	the description of each given environment: {env_description}. A unit operation you
1231	can perform is called action in a given environment. For each environment, you are given a
1232	limited action space as function calls:
1233	{action_descriptions}
1234	You may receive a screenshot of the current system. You may receive a screenshot of a
1235	smartphone app. The interactive UI elements on the screenshot are labeled with numeric tags
1236	starting from 1.
1237	the next action, then use a provided action in each step to achieve the task. You should state
1238	what action to take and what the parameters should be. Your answer MUST be a least one
1239	function call. You SHOULD NEVER ask me to do anything for you. Always do them by
1240	vourself using function calls.
1241	

1242	Prompt
1243	Tompt
1244	You are a helpful assistant. Now you have to do a task as described below:
1245	**{task_description}**
1246	You should never forget this task and always perform actions to achieve this task. And this
1247	is the description of each given environment: {env_description}. You will receive
1248	screenshots of the environments. The interactive UI elements on the screenshot are labeled
1249	with numeric tags starting from 1.
1250	A unit operation you can perform is called Action. You have a limited action space as function
1251	calls: {action_descriptions}. You should generate JSON code blocks to execute the actions. Each code block MUST contains only one ison object i.e. one action. You can output
1252	multiple code blocks to execute multiple actions in a single step. You must follow the ISON
1253	format below to output the action.
1254	{"name": "action_name", "arguments": {"arg1": "value1",
1255	"arg2": "value2"}}
1256	or if not arguments needed:
1257	{"name": "action_name", "arguments": {}}
1258	You MUST use exactly the same "action_name" as I gave to you in the action space. You
1259	SHOULDN'T add any comments in the code blocks.
1260	In each step, You MUST explain what do you see from the current observation and the plan of
1261	the next action, then use a provided action in each step to achieve the task. You should state
1262	one code block. You SHOLILD NEVER ask me to do anything for you. Always do them by
1263	vourself
1264	

### 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.

### B.3.3 MULTI-AGENT BY ENVIRONMENT

### 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.

1299

1300

1301 1302 1303

1304

1305

1306

1309

1310

1311

1312

1313

1315

1316

1317

1318

1319 1320 1321

1322 1323

1327

#### 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

1324This section further discusses our experimental results in detail. Section C.1 categorizes the results into<br/>three types of tasks: Ubuntu, Android, and cross-platform, and provides further analysis. Section C.3<br/>examines three specific tasks and analyzes the performance of different agent settings on each.

1328 C 1 RESI

C.1 RESULT BY PLATFORMS

Table 7, 8 and 9 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-40, 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.

1338 Cross-platform tasks necessitate functionality across different operating systems or platforms, demand 1339 a broader capability range and more sophisticated agent coordination. The importance of CR is 1340 especially critical in such environments, where it serves as a more reliable metric for distinguishing 1341 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 1342 capability, thereby better reflecting its robustness and adaptability to complex requirements. On 1343 cross-platform tasks, GPT-4 Turbo (Single) exhibits a CR of 52.61%, which indicates that even 1344 though SR might be lower, the agent covers a significant portion of task objectives before termination. 1345

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-40 with JSON mode

Table 7: Evaluation results on Ubuntu tasks.								
Agent system			M	Termination Reason				
Model	Structure	SR(%) ↑	$  CR(\%) \uparrow$	<b>EE(%)</b> ↑	<b>CE(%)</b> ↑	FC(%)	RSL(%)	IA(%)
GPT-40	Single	9.59	30.82	3.22	$4.87 \times 10^{-4}$	6.85	58.90	24.66
GPT-40	By Func	9.59	24.20	2.72	$4.30 \times 10^{-4}$	5.48	63.01	21.92
GPT-40	By Env	10.96	22.88	2.74	$2.29 \times 10^{-4}$	5.48	43.84	39.73
GPT-4 Turbo	Single	10.96	31.09	4.08	<b>5.57</b> $\times 10^{-4}$	2.74	65.75	20.55
GPT-4 Turbo	By Func	12.33	28.95	3.70	$4.18 \times 10^{-4}$	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	$0.63 \times 10^{-4}$	8.22	58.90	32.88
CLAUDE 3 OPUS	By Func	0.00	4.93	0.46	$0.47 \times 10^{-4}$	27.40	34.25	38.36
GPT-40 w/o FC	Single	10.96	22.58	2.30	$4.49 \times 10^{-4}$	5.48	54.79	28.77
PIXTRAL-12B	Single	0.00	2.97	0.22	$0.24 \times 10^{-4}$	1.37	80.82	17.81
LLAVA-OV-72B	Single	0.00	3.31	0.20	$0.35 \times 10^{-4}$	17.81	64.38	17.81

1350

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.

1376 1377

### 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 1382 to a high False Completion Rate—where the agent incorrectly assumes that the task is complete. As 1383 observed in failure cases (e.g., the Cross-platform Task case study in Appendix C.3), Sub-agents 1384 often misinterpret the Main agent's instructions. Despite being required to perform a final action, the 1385 instructions lead Sub-agents to prematurely conclude that the task is complete, resulting in incorrect 1386 "complete" actions. While this issue also occurs in Multi-Env, it happens less frequently. By analysing 1387 the communication logs, we believe this is due to information loss during inter-agent communication. 1388 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 1389 in LLM communication, is less suited for inter-agent communication, leading to information loss 1390 during compression and interpretation, which weakens the performance of multi-agent structures. 1391

1392 Next, comparing in Invalid Action Rate, we observe that in single-platform tasks, both Multi-Env and 1393 Multi-Func suffer from similar inter-agent communication issues, as indicated by their high Invalid 1394 Action rates. However, in cross-platform tasks (Table 9), the Single agent's Invalid Action rate is 1395 significantly higher than that of the Multi-agent by environment structures on GPT-40 model. Crossplatform tasks require frequent environment changes with varying action spaces, and if the model's 1396 performance output is inadequate, it often generates correct actions in the wrong environment, invalid 1397 actions in the correct environment, or correct actions in correct environment but in the wrong format. 1398 This phenomenon highlights the limitations of current general-purpose LLMs, where multi-agent 1399 structures can be advantageous. By assigning each agent a specific responsibility and a limited action 1400 space, multi-agent structures can mitigate these issues. 1401

 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

Table 8: Evaluation results on Android tasks.									
Agent system			M	Termination Reason					
Model	Structure	SR(%) ↑	$  CR(\%) \uparrow$	<b>EE(%)</b> ↑	<b>CE(%)</b> ↑	FC(%)	RSL(%)	IA(%)	
GPT-40	Single	24.14	47.91	5.84	$7.17 \times 10^{-4}$	13.79	58.62	3.45	
GPT-40	By Func	24.14	48.74	6.83	<b>9.19</b> ×10 <sup>-4</sup>	24.14	37.93	13.79	
GPT-40	By Env	27.59	53.34	6.99	$4.58 \times 10^{-4}$	13.79	44.83	13.79	
GPT-4 Turbo	Single	6.90	27.08	2.60	$2.87 \times 10^{-4}$	20.69	62.07	10.34	
GPT-4 Turbo	By Func	20.69	37.01	5.00	$5.92 \times 10^{-4}$	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	$5.37 \times 10^{-4}$	20.69	55.17	10.34	
CLAUDE 3 OPUS	By Func	13.79	44.02	4.75	$5.35 \times 10^{-4}$	48.28	31.03	6.90	
GPT-40 w/o FC	Single	10.34	14.29	1.72	$2.94 \times 10^{-4}$	3.45	6.90	79.31	
PIXTRAL-12B	Single	3.45	24.17	2.16	$2.72 \times 10^{-4}$	0.00	65.52	31.03	
LLAVA-OV-72B	Single	3.45	13.51	1.36	$3.00 \times 10^{-4}$	3.45	93.10	0.00	

Table 9: Evaluation results on cross-platform tasks.

Agent system			M	Termination Reason				
Model	Structure	<b>SR(%)</b> ↑	$  CR(\%) \uparrow$	<b>EE(%)</b> ↑	<b>CE(%)</b> ↑	FC(%)	RSL(%)	IA(%)
GPT-40	Single	16.67	51.24	5.21	<b>3.98</b> $\times 10^{-4}$	5.56	38.89	38.89
GPT-40	By Func	22.22	50.00	4.15	$3.13 \times 10^{-4}$	11.11	44.44	22.22
GPT-40	By Env	5.56	43.54	3.22	$1.60 \times 10^{-4}$	11.11	72.22	11.11
GPT-4 Turbo	Single	5.56	52.61	4.60	$2.89 \times 10^{-4}$	11.11	66.67	16.67
GPT-4 Turbo	By Func	5.56	46.17	4.06	$2.67 \times 10^{-4}$	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	$1.24 \times 10^{-4}$	0.00	55.56	44.44
CLAUDE 3 OPUS	By Func	0.00	18.96	1.93	$1.20 \times 10^{-4}$	0.00	38.89	61.11
GPT-40 w/o FC	Single	0.00	39.11	3.51	$3.28 \times 10^{-4}$	5.56	50.00	44.44
PIXTRAL-12B	Single	0.00	12.35	0.62	$0.44 \times 10^{-4}$	0.00	72.22	27.78
LLAVA-OV-72B	Single	0.00	9.07	0.48	$0.53 \times 10^{-4}$	5.56	66.67	27.78

<sup>1438</sup> 

1404

1439

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.

#### 1443 1444 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.

1449 C.3.1 CROSS-PLATFORM TASK 1450

1451Task: Open the "Tasks" app on an Android device, check the first incomplete task, and then<br/>execute it as described. The first task, found incomplete in the "Tasks" app, involves switching the<br/>system to dark mode in Ubuntu via the "Settings" application.

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

1458 In this task, GPT-40 (single agent), GPT-4 Turbo (single agent), and GPT-4 Turbo (multi-agent by 1459 functionality) all successfully complete the task using the minimal steps necessary to locate and exe-1460 cute the task, demonstrating their efficiency in managing multiple environments simultaneously. On 1461 the other hand, both GPT-40 (multi-agent by functionality) and GPT-40 (multi-agent by environment) 1462 also perform commendably, completing the task up until the final step. However, after incorrectly performing the last step, they both erroneously conclude the task is completed and exit. This indicates 1463 a communication breakdown, where the sub-agents misinterpret the instructions from the main agent. 1464 The remaining four agents fail to complete the task. Agents equipped with the Gemini model do 1465 not even manage to open the "Tasks" app within the allocated step limit, whereas agents with the 1466 Claude model quickly open the "Tasks" app to complete the first step but fail at the task execution. 1467 The performance disparity between single-agent and multi-agent configurations in both the Gemini 1468 and Claude models highlights the variance in capability across different models and devices. 1469

1470

### 1471 C.3.2 UBUNTU TASK

## 1472<br/>1473Task: Create a new directory "/home/crab/assets\_copy" and copy all files with the specified<br/>"txt" extension from "/home/crab/assets" to the directory "/home/crab/assets\_copy".

This task can be approached through multiple methods. An agent may opt for a straightforward 1475 strategy first using the search\_application command to find the Terminal, then using Linux 1476 commands to create the directory and copy the necessary files. Alternatively, the agent could employ 1477 a GUI-based approach, manually creating the folder and selecting files through actions like click 1478 and right\_click. We evaluate various agent systems in a single-agent setting for this task. As 1479 illustrated in Table 10-13, both GPT-40 and GPT-4 Turbo from OpenAI successfully interpret the task 1480 instructions and employ a simpler solution using Terminal commands. These agents also demonstrate 1481 superior capability in understanding the UI, selecting the correct commands, and accurately using the 1482 Terminal application to fulfill the task requirements.

1483 Conversely, the Gemini and Claude agents, despite attempting to solve the task with Terminal, 1484 ultimately fail in different ways. Both agents struggle with precise clicking and selecting the correct 1485 icons for the intended actions, even though they share the same visual prompting mechanism as 1486 GPT-40 and GPT-4 Turbo. For instance, the Claude agent mistakenly opens the Ubuntu Desktop 1487 Guide instead of the Terminal and continues executing commands in the wrong application without 1488 realizing the error. The Gemini agent, on the other hand, unexpectedly opens the Firefox browser before correctly navigating to the Terminal but still interacts incorrectly with unrelated applications 1489 and icons. Unlike Claude, Gemini does not type in commands in the wrong applications but persists 1490 in exploring alternative methods using the Files application's UI. Despite taking significantly more 1491 steps than the GPT-40 and GPT-4 Turbo agents, neither the Claude nor the Gemini agents achieve the 1492 task's goal. 1493

1494

### 1495 C.3.3 ANDROID TASK

## 1496<br/>1497<br/>1498Task: In Android, using the "Contacts" app, find the email of the contact named John Lauphin,<br/>then using the "Gmail" app, send an email to that contact with the subject "Hello John."

This task consists of sub-tasks across two different applications. Agents must sequentially open the two apps, retrieve the email address from the first app, and use it in the second app to send an email. 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 reports the performance of agents in a multi-agent setting for this challenging task. Following is the details of agents in operating the task.

1505

**GPT-40 multi-agent by functionality** In steps 1-11, the agent tries to open the Contacts app but
 mistakenly opens Google Assistant multiple times. In steps 12-14, the agent successfully enters the
 Contacts app and finds the contact information. The agent then returns to the home page, and the
 process is terminated due to the limitation of operation steps.

- 1510
- **GPT-4 Turbo multi-agent by functionality** In steps 1-2, the agent tries to open the Contacts app but mistakenly opens Google Messages. In steps 3-5, the agent opens the Contacts app and obtains the

corresponding information. In steps 6-14, the agent repeatedly opens Google Chrome and Messages apps, failing to find the Gmail app as planned.

1515 Gemini 1.5 Pro multi-agent by functionality In steps 1-2, the agent finds the Contacts app and
 1516 enters it. However, the agent misunderstands the instruction, gets lost in creating a new contact with
 1517 the given name, and cannot obtain the corresponding information.

1519 Claude 3 Opus multi-agent by functionality In steps 1-7, the agent tries to open the Contacts app but mistakenly opens Google Messages multiple times. In steps 7-11, the agent tries to open the Contacts app but mistakenly opens Google Assistant. In steps 12-14, the agent successfully enters the Contacts app and finds the contact information. The agent then returns to the home page, plans to open the Gmail app, and the process is terminated due to the limitation of operation steps.

**GPT-40 multi-agent by environment** In steps 1-7, the agent plans to open the Contacts app, but the operation fails due to an error in opening the app drawer, which prevents the agent from finding and tapping the Contacts app. In steps 8-11, the agent successfully enters the Contacts app and obtains the information. In steps 12-14, the agent opens the Gmail app, navigates to the sending page, and tries to input the retrieved email address as the recipient.

**Analysis** For the agents which are organized by functionality, Gemini 1.5 Pro struggles to complete the first operation. Although it recognizes and opens the Contacts app as instructed, it fails to proceed further. In contrast, Claude 3 Opus and GPT-40 successfully obtain the necessary information. In the initial phase, the multiple agents agree that opening the Contacts app is the first step. However, they often fail to find the correct position to tap, frequently opening incorrect apps such as Google Assistant and Messages. Once the agents do open the correct app, they usually find the email address of the contact quickly. Even when agents plan to go back home and open the Gmail app to send the message, due to the limitation of operations, the system ended. As shown in steps 3-5, GPT-4 Turbo quickly finishes the corresponding task after opening the correct app. However, similar to GPT-40, GPT4-Turbo agents get stuck as they can not open the correct apps in the following steps. Besides, GPT-40 (multi-agent by environment) overcomes the issue encountered by GPT-40 (multi-agent by functionality). Even affected by not being able to access the app drawer, the system could still find and copy the corresponding information and change to the Gmail app for further operations.

Table 10: **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".











Table 11: **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".







Table 12: **Ubuntu task case with GPT-40** (**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".







Table 13: **Ubuntu task case with GPT-4 Turbo (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".





