

000 SCALING SYNTHETIC TASK GENERATION FOR AGENTS 001 002 VIA EXPLORATION 003 004

005 **Anonymous authors**

006 Paper under double-blind review

007 008 ABSTRACT 009

011 Post-Training Multimodal Large Language Models (MLLMs) to build interactive
012 agents holds promise across domains such as computer-use, web navigation, and
013 robotics. A key challenge in scaling such post-training is lack of high-quality
014 downstream agentic task datasets with tasks that are diverse, feasible, and verifiable.
015 Existing approaches for task generation rely heavily on human annotation or
016 prompting MLLM with limited downstream environment information, which is
017 either costly or poorly scalable as it yield tasks with limited coverage. To remedy
018 this, we present AUTOPLAY, a scalable pipeline for task generation that explicitly
019 explores interactive environments to discover *possible interactions* and *current state*
020 information to synthesize environment-grounded tasks. AUTOPLAY operates in
021 two stages: (i) an exploration phase, where an MLLM explorer agent systematically
022 uncovers novel environment states and functionalities, and (ii) a task generation
023 phase, where a task generator leverages exploration trajectories and a set of task
024 guideline prompts as context to synthesize diverse, executable, and verifiable tasks.
025 We show AUTOPLAY generates 20k tasks across 20 Android applications and 10k
026 tasks across 13 applications Ubuntu applications to train mobile-use and computer-
027 use agents. AUTOPLAY generated tasks enable large-scale task demonstration
028 synthesis without human annotation by employing an MLLM task executor and
029 verifier. This data enables training MLLM-based UI agents that improve success
030 rates up to 20.0% on mobile-use and 10.9% on computer-use scenarios. In addition,
031 AUTOPLAY generated tasks combined with MLLM verifier-based rewards enable
032 scaling reinforcement learning training of UI agents, leading to an additional 5.7%
033 gain. coverage. These results establish AUTOPLAY as a scalable approach for
034 post-training capable MLLM agents reducing reliance on human annotation.

035 1 INTRODUCTION 036

037 Multimodal Large Language Models (MLLMs) are a promising foundation for building agents
038 across a wide range of downstream domains, including computer use (Qin et al., 2025; Agashe
039 et al., 2025; OpenAI, 2025), web navigation (Zhou et al., 2023; Yao et al., 2024; 2022), video
040 games (Fan et al., 2022; Wang et al., 2023), and robotics (Driess et al., 2023; Black et al., 2024;
041 Kim et al., 2024). Owing to their broad knowledge and reasoning capabilities, these models are
042 well-suited to understanding and planning the execution of open-ended tasks across these diverse
043 application areas. A central challenge in training such agents is the scarcity of downstream interactive
044 agentic data. Such interactive data consists of two components: (1) Diverse Tasks: a sufficiently
045 broad set of tasks or queries that cover real-world use cases of such agents, (2) Task demonstrations:
046 corresponding task execution trajectory. For most of these domains, such data is not readily available
047 at web-scale, as much of the relevant interaction data resides on closed systems such as personal
048 devices and commercial hardware. Consequently, post-training efforts for building agentic MLLMs
049 for such domains have relied heavily on human annotation to source a large pool of diverse tasks
050 and corresponding demonstrations (Li et al., 2024; Wang et al., 2025; Qin et al., 2025). However,
051 this approach is prohibitively expensive and scales poorly. We argue, the fundamental bottleneck to
052 enable scalable post-training of agentic MLLMs is lack of large high-quality task definition datasets
053 with tasks that are diverse, feasible to execute, verifiable, and aligned with real-world use cases. With
access to such task datasets, MLLMs could be post-trained in a scalable manner with synthetically
generated demonstrations using supervised finetuning (SFT) or via Reinforcement Learning (RL)

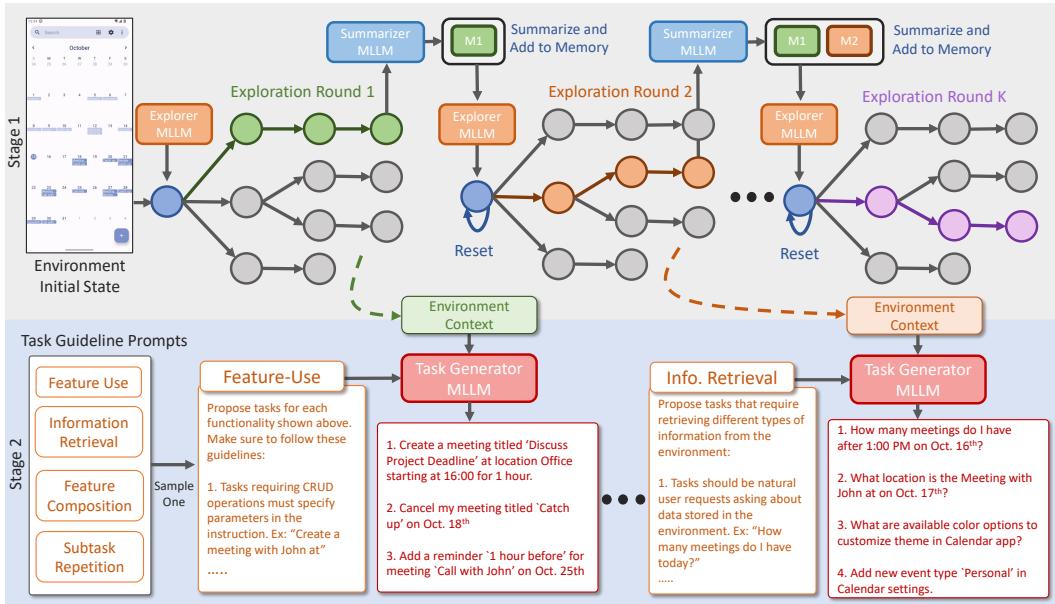


Figure 1. AUTOPLAY generates large-scale, diverse and verifiable tasks for scaling supervision for MLLM agents. In stage 1 (top), AUTOPLAY covers the environment states through a MLLM exploration policy that tracks seen states via a memory module. Next, stage 2 (bottom) uses these exploratory trajectories and task guideline prompts as context for proposing tasks. The guidelines help enforce task diversity and the exploration trajectories uncover environment features and content relevant for proposing tasks.

without any human annotation. This highlights the need for an automatic and reliable pipeline to synthesize tasks at scale.

For synthesized task datasets to be useful, they must provide broad *coverage* of the target environment, ensuring a diverse and representative training set. They must also be *feasible*, meaning that agents can realistically execute them to generate useful demonstrations. Finally, they should be *verifiable* to filter high-quality trajectories for SFT or to build a reward model for RL. However, building an automatic pipeline that generates tasks with these properties is challenging as it requires explicit knowledge of current state information and what interactions are feasible in an environment – knowledge that can only be obtained through direct interaction with the environment.

In an attempt to address this challenge, prior works rely on limited domain knowledge of an MLLM to generate the tasks (Trabucco et al., 2025; Zhou et al., 2025), this typically yields generic tasks with limited coverage and offers no guarantees of feasibility or verifiability as current MLLMs lack explicit grounding in both environment dynamics and current state information. To address these shortcomings, we propose an approach that actively explores the agent’s environment in an exhaustive manner, collecting relevant information that can be used to ground MLLMs in the environment for scalable task generation. In this work, we focus on UI agents (Rawles et al., 2024; Xie et al., 2025) that observe the current screenshot and interact with UI elements as a human user would. This setting is both broad, capturing the full spectrum of everyday device tasks and practical since there are programmatic environments available for UI interactions (Rawles et al., 2024; Xie et al., 2024).

We introduce AUTOPLAY an approach for scalable task generation with a focus on task coverage, feasibility, and verifiability. Our method, depicted in Figure 1, incorporates two phases of exploration and task generation. First, in *environment exploration* phase, an MLLM explorer agent equipped with memory is prompted to exhaustively explore an increasing number of novel environment states (top of Fig. 1). Such exploratory trajectories are intended to discover the accessible functionalities and content of the environment. Next, in *task generation* phase, a task generator MLLM uses exploration trajectories as environment context to produce diverse environment-grounded tasks based on a set of task guideline prompts which describe desired task properties (bottom of Fig. 1). For instance, a task guideline prompt for Feature-Use tasks would encourage generation of tasks that require doing diverse create, edit, of delete operations on entities in the environment. We present an example of the exploration trajectory collected by AUTOPLAY and the corresponding tasks synthesized using task generator grounded in the state of the environment in Fig. 2.

108 We use AUTOPLAY to scale task generation for mobile and computer UI agents. AUTOPLAY generates
 109 20k tasks across 20 apps in an Android platform and 10k tasks across 13 apps in an Ubuntu platform.
 110 We then synthesize demonstrations for the AUTOPLAY-generated tasks using an MLLM executor and
 111 verify them with an MLLM verifier, without relying on privileged environment information or human
 112 annotation. These demonstrations are used to finetune an MLLM agent with SFT. Additionally,
 113 AUTOPLAY generated tasks also enable training the MLLM agent with RL, using the MLLM verifier
 114 as a reward model. Through these experiments we demonstrate AUTOPLAY enables post-training
 115 MLLM agents using both SFT and RL in a scalable manner without relying on any human annotation.

116 We show that the AUTOPLAY pipeline is able to train effective MLLM UI agents. In mobile-use,
 117 AUTOPLAY boosts success rate performance by 13%–20% over the base model across a range of
 118 model sizes. In computer-use, our method improves base model success rate performance by up
 119 to 10.9%. Beyond SFT, incorporating the generated tasks with the verifier for RL training yields
 120 an additional 5.7% improvement in task success rate. Together, these results demonstrate that
 121 AUTOPLAY-generated tasks and trajectories lead to consistent and significant performance gains
 122 across model sizes and training paradigms. We also find that the AUTOPLAY task generation pipeline
 123 significantly outperforms prior approaches for synthetically generating tasks in the mobile-use
 124 domain (Trabucco et al., 2025; Zhou et al., 2025; Pahuja et al., 2025; Xie et al., 2025) which involve
 125 either generating tasks with limited environment information or methods that interleave task proposal
 126 with execution. We find this improvement is a result of AUTOPLAY generating tasks with higher
 127 diversity, coverage, and feasibility for task execution than prior methods, ultimately yielding datasets
 128 that better support the training of capable UI agents.

2 RELATED WORK

130 **Zero-Shot Methods for Agentic Tasks.** Modular agentic pipelines (Agashe et al., 2025; Koh et al.,
 131 2024) that leverage MLLMs as high-level task planners, combined with diverse tools (*e.g.* low-level
 132 action controllers, memory, and tool use), have emerged as an effective approach for building capable
 133 UI agents across various domains (Xie et al., 2024; Rawles et al., 2024; Yao et al., 2024; Deng
 134 et al., 2023; Zhou et al., 2023). These pipelines aim to decompose the skills required for complex
 135 tasks—such as mapping high-level actions to environment-specific low-level controls, maintaining
 136 interaction history, verifying execution outcomes, and invoking tools—into specialized modules. Such
 137 modules can then be flexibly composed, enabling scalable and adaptable modular agentic pipelines.
 138 In this work, we leverage these agentic pipelines as synthetic data generation modules that enable
 139 us to collect UI-interaction demonstrations for AUTOPLAY generated tasks to bootstrap training
 140 end-to-end UI agents that take environment observation as input and directly output actions.

141 **Synthetic Task and Trajectory Generation.** Synthetic data generation using agentic pipelines has
 142 emerged as a promising approach to unlock internet-scale data for training UI agents. For instance,
 143 methods like PAE Zhou et al. (2025) leverage LLMs conditioned on limited information about the
 144 environment (*i.e.* manually written textual descriptions) to propose tasks to synthesize task execution
 145 trajectories. In contrast, AUTOPLAY explicitly explores the environment to gather richer context
 146 about the environment without relying on human annotated textual descriptions. Another line of
 147 research focuses on iteratively proposing and executing tasks. Xie et al. (2025); Pahuja et al. (2025)
 148 use an initial screenshot from the environment and propose a short-horizon subtask, and then execute
 149 this subtask, repeating this process and using the summary of previously executed sub-task as context
 150 to propose the next subtask. This chain of subtasks is then summarized into a single long-horizon task
 151 with hindsight relabeling. Similarly, Trabucco et al. (2025) iteratively proposes and executes subtasks,
 152 but only for a single iteration without chaining subtasks. Like AUTOPLAY these methods ground tasks
 153 in environment interactions, yet they are limited to each trajectory directly mapping to an instruction
 154 via hindsight relabeling. AUTOPLAY can propose multiple tasks from a single trajectory based on the
 155 uncovered possible interactions and sampled task guideline prompt. Furthermore, these prior works
 156 require each iterative subtask to start from the previous subtask, which progressively constrains the
 157 space of exploratory trajectories. On the other hand, methods like (Murty et al., 2024; 2025) start
 158 by using an LLM as a unconstrained exploration policy, next a trajectory labeler summarizes the
 159 exploration trajectories at regular intervals to a task instruction. Next, a verifier verifies whether
 160 the trajectory completes the proposed task instruction or not. If it does not, a instruction following
 161 policy attempts to execute the proposed task instruction. In contrast to our method, (Murty et al.,

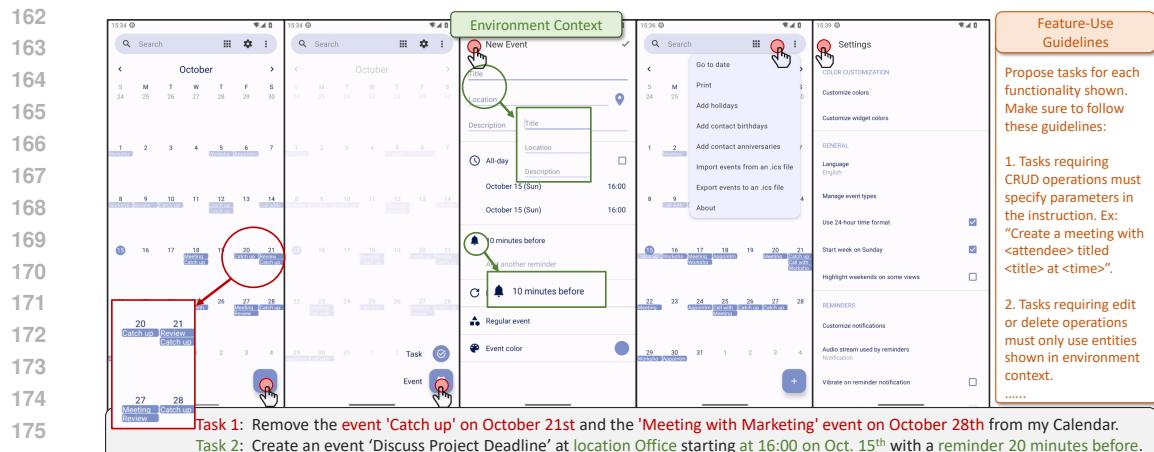


Figure 2. Example of task generation based on an environment context, represented as a set of screenshots and interactions, and a set of task guidelines. AUTOPLAY uses presences of events in the calendar together with guidance of using entities in the context, such as names and dates, to produce Task 1. Similarly, showing the event creation form in the context, coupled with guidance to use these fields as task parameters, results in Task 2.

2024; 2025) does not use an exploration policy with memory and explicit task guidelines to maximize coverage of the environment or maximizing task diversity.

3 AUTOPLAY

3.1 PRELIMINARIES

We define AUTOPLAY in the context of multi-step decision making domains that be expressed as Partially Observable Markov Decision Processes (POMDP) (Puterman, 2014). A POMDP can be defined as a tuple $(\mathcal{S}, \mathcal{O}, \mathcal{A}, P, R, \rho_0)$ for underlying state space \mathcal{S} , observation space \mathcal{O} , action space \mathcal{A} , transition function $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$, reward function R and initial state distribution ρ_0 . We consider the extension of including a goal distribution \mathcal{G} and the case where the reward is formulated as $R(s, g)$ for $s \in \mathcal{S}$ and $g \in \mathcal{G}$. We seek to learn a goal-conditioned policy $\pi(a_t|o_t, g)$ mapping from observation o_t at timestep t and goal g to an action a_t to achieve the goal state s_g . We train a goal-conditioned policy $\pi(a|o, g)$ in a POMDP with a dataset of tasks \mathcal{D} . Using the above formalism, dataset of task \mathcal{D} consists of tuples (g, s_0) where g is a goal-specification in natural language, and s_0 is the initial state of the environment. Similarly, we define trajectory as a sequence $\tau = (o_0, a_0, \dots, o_T)$ where the outcome of achieving goal g is verified by a reward model. In our case, the reward model only relies on the trajectory $R(\tau, g)$ without privileged access to the environment state.

The goal of AUTOPLAY is to automatically generate large dataset of tasks \mathcal{D} for a specific domain by actively interacting with POMDP to gather state information, semantics, and understanding dynamics. With access to task dataset \mathcal{D} one can either attempt to generate training trajectories using a data collection policy for supervised finetuning (SFT), or use the reward model R in a multi-task reinforcement learning (RL) setup to train a goal-conditioned policy $\pi(a|o, g)$. We instantiate AUTOPLAY in the UI agent domain for mobile-use and computer-use environments, e.g. a Calendar app on a mobile/desktop device. In this domain a state observation \mathcal{O} is the partially observable device screenshot, an action is a UI interaction such as ‘click’, ‘type’, ‘scroll’, etc., and the transition operator is defined by transition dynamics of the UI application. A task is defined by a goal like $g = \text{"Remove all the events on my calendar for next Tuesday"}$ and initial state s_0 starts with the agent on the homepage of the application with data populated in the Calendar application.

3.2 AUTOPLAY TASK GENERATOR

AUTOPLAY task generator operates in two stages: First, the POMDP is explored in a goal-agnostic manner to maximize coverage over novel states of the environment using an explorer agent. Second, the task dataset \mathcal{D} is produced using the information of the explored states combined with domain specific task guidelines.

216	Algorithm 1 AUTOPLAY	
217	Stage 1: Environment Exploration	
218	Parameters:	Stage 2: Task Generation
219	N # of apps	Parameters:
220	M : # of exploration turns	\mathcal{P} : task guidelines
221	$\mathcal{E} = \emptyset$	K # of tasks per guideline and context
222	for $j = 1 \dots N$ do	$\mathcal{D} = \emptyset$
223	Sample $s_0 \sim \mathcal{S}$, Initialize context $M = \emptyset$	for $s = 1 \dots S$ do
224	for $k = 1 \dots M$ do	Sample $p \sim \mathcal{P}$, $\tau \sim \mathcal{E}$
225	Sample trajectory τ using MLLM(M , explorer_prompt)	Sample (g_1, \dots, g_K) using
226	Summarize as $m = \text{MLLM}(\tau, \text{summary_prompt})$	MLLM(task_generator_prompt, p, τ, \cdot)
227	$\mathcal{E} = \mathcal{E} \cup \{\tau\}$, $M = M \cup \{m\}$	$\mathcal{D} = \mathcal{D} \cup \{(g_1, s_1), \dots, (g_K, s_1)\}$
228	end for	where s_1 is the first state in τ
229	end for	end for
230		return \mathcal{D}
231		

228 **Stage 1: Environment Exploration.** In this stage, a goal-agnostic explorer agent is employed to
 229 exhaustively explore the environment to maximize coverage over novel states and functionalities
 230 of the environment. The explorer agent is implemented using an MLLM agent equipped with
 231 an explicit memory of past interactions. At each timestep, the MLLM agent is provided current
 232 environment observation, past interaction memory, and prompted to select diverse actions, aiming
 233 to cover the full range of possible interactions within an environment. The explorer agent interacts
 234 with the environment for K steps. This process produces a exploration trajectory, denoted as
 235 $\tau = (o_1, a_1, \dots, o_K, a_K)$, which we refer to as the *environment context*.

236 To ensure the environment is explored as exhaustively as possible, we repeat the exploration process
 237 for each environment M times. To encourage the explorer to cover novel states we provide explorer
 238 access to prior exploration turns $\{\tau_1, \dots, \tau_{i-1}\}$ *i.e. episodic memory*. As MLLMs have finite
 239 context window, directly conditioning on entire high-dimensional past trajectories is infeasible.
 240 Therefore, for each exploration turn i , we summarize the prior exploration turns $\{\tau_1, \dots, \tau_{i-1}\}$ to
 241 $\{m_1, \dots, m_{i-1}\}$ where m is a concise and comprehensive text representation describing exploration
 242 trajectory τ generated using a summarizer MLLM by $m = \text{MLLM}(\text{summary_prompt}, \tau)$. At the
 243 end of stage 1, we obtain a full *environment context* $\mathcal{E} = \{\tau_1, \dots, \tau_M\}$.

244 **Stage 2: Exploration Conditioned Task Generation** In this stage, each environment context
 245 from $\tau \sim \mathcal{E}$ is used to define a set of N tasks. We would like to highlight that, the environment
 246 context τ does not represent a successful execution of a task, rather it serves as context of what
 247 interactions are feasible and current state of the environment which can be used as hints to synthesize
 248 tasks. More concretely, the environment context τ helps uncover underlying information of the
 249 POMDP transition function \mathcal{T} and state space \mathcal{S} . We use a MLLM as a task generator that uses
 250 τ to generate environment-grounded tasks. Since there are many possible ways to derive task
 251 instructions from a single environment context, we find it useful to provide domain-specific prompts
 252 that specify guidelines for what constitutes a good task to our task generator MLLM. These prompts,
 253 $\mathcal{P} = \{p_1, \dots, p_N\}$, referred to as *task guidelines*, are tailored to the target domain. For the UI domain,
 254 they are illustrated in Fig. 2. For example, one guideline, called Feature-Use, encourages creation of
 255 tasks that require doing diverse create, edit, of delete operations on entities in the environment shown
 256 in the corresponding environment context.

257 To produce the task dataset \mathcal{D} , we append the task generator prompt, with trajectory data and
 258 task guideline prompt. We iterate over all combinations of trajectories in environment con-
 259 text $\tau \in \mathcal{E}$ and task guideline prompt $p \in \mathcal{P}$ to sample a sequence of goals $(g_1, \dots, g_N) \sim$
 260 $\text{MLLM}(\text{task_generator_prompt}, p, \tau)$. Each goal, when paired with the initial state of τ , results in
 261 the task dataset. An example of this process for the UI domain is shown in Fig 2.

262 The above approach is summarized in Algorithm 1. The prompts used for the explorer, summarizer,
 263 and task generator MLLMs are described in Appendix D.2. The explorer agent follows the same
 264 architecture as the task executor agent detailed in Section 3.3.

265 3.3 TRAINING DATA GENERATION AND VERIFICATION

266 In the following, we describe the task executor and verifier used in the UI domain.

267 **Task Executor Agent.** We use the tasks generated by AUTOPLAY to produce successful trajectories
 268 via an MLLM based executor agent. Our executor agent builds on the system proposed in [Yang](#)

270 [et al. \(2025\)](#) and consists of an MLLM planner, a reflection tool, and a grounding model. Given the
 271 task instruction, the current screenshot, previously executed actions, and a reflection trace of the
 272 last action generated by reflection tool, the MLLM planner generates a high-level action in natural
 273 language. If this high-level action is coordinate-based (e.g., a click), the grounding model translates it
 274 into a pixel-level coordinate for execution. Non-coordinate-based actions are executed directly. The
 275 reflection tool supplements this process by taking as input the previous action, the prior observation,
 276 and the current observation to generate a reflection trace. This trace describes the effect of the action
 277 on the environment and whether it was executed successfully, and it is provided as additional context
 278 to the high-level policy at the next timestep.

279 For both mobile-use and computer-use domains, we use GPT-4o ([OpenAI et al., 2024](#)) as both the
 280 planner and reflection model. We use UI-TARS-1.5 7B ([Qin et al., 2025](#)) as the grounding model
 281 for mobile and GTA1-7B ([Yang et al., 2025](#)) as the grounding model for computer-use. In the
 282 computer-use setting, we additionally supply the high-level policy with one heuristic action, to enable
 283 interaction with complex UI elements (like spreadsheets, detailed in [Appendix G](#)), which improves
 284 data collection success rates. These heuristic actions are only used by the expert and are not available
 285 to the trained AUTOPLAY policy.

286 **Task Trajectory Verifier.** To determine whether AUTOPLAY-generated tasks are executed success-
 287 fully, we employ a task verifier that evaluates trajectories. The verifier is an MLLM that takes as
 288 input the task instruction and the executed trajectory (represented as interleaved images and actions).
 289 It operates in three stages: (i) summarizing what the agent is doing in the trajectory, (ii) producing
 290 a Chain-of-Thought ([Wei et al., 2023](#)) reasoning for whether the task has been completed, and (iii)
 291 issuing a final judgment of “success” or “failure”. We use GPT-4o as the task verifier for trajectories
 292 collected using the task executor agent. Full details of the verifier and prompts are in [Appendix D.4](#).

293 **SFT Data Generation.** To generate SFT data for mobile-use domain, we use 20 mobile apps
 294 from [Rawles et al. \(2024\)](#). Similarly, for computer-use domain, we use 13 desktop apps from [Xie et al.](#)
 295 (2024) on Ubuntu devices. We run the environment exploration for 3 turns to build the environment
 296 context and we define 4 task guideline prompts. Using these, we generate 50 tasks for each app, task
 297 guideline prompt, exploration context combination. We specify the precise prompts used for the
 298 task guidelines in [Appendix D.1](#) for each platform. We are able to generate $\sim 20k$ tasks for Android
 299 platform and $\sim 10k$ for Ubuntu platform. After execution of these tasks and subsequent verification
 300 with the verifier we obtain $\sim 8k$ and $\sim 3.5k$ successful trajectories respectively to use for SFT.

301 **Base Model.** For our experiments, we use Qwen2.5-VL Instruct ([Bai et al., 2025](#)) 3B-72B MLLMs
 302 as the base model and finetune them using SFT and RL on our dataset to build a end-to-end UI agents.
 303 At each timestep, the policy takes as input: current image observation, history of past actions, the
 304 task instruction and outputs low-level actions. [Appendix G](#) describes the action space per-domain.

3.4 REINFORCEMENT LEARNING ON SYNTHETIC TASKS

307 We also use the tasks generated by AUTOPLAY in conjunction with the task verifier to scale rein-
 308 force learning (RL) training of mobile-use agents. For each RL training environment worker, we
 309 sample a random task from the AUTOPLAY task dataset $(g, s_0) \sim \mathcal{D}$. We then initialize the simulator
 310 state to s_0 and task our agent to successfully complete the instruction. At the end of the rollout, we
 311 score the trajectory with the task verifier as either a success or failure, assigning a reward of 1 or
 312 0 to the trajectory respectively. We use the Qwen2.5-VL-32B-Instruct [Bai et al. \(2025\)](#) MLLM to
 313 instantiate the task verifier using the same verifier setup as in [Sec. 3.3](#). By leveraging the AUTOPLAY
 314 task generator and an MLLM task verifier we unlock RL training with verifier feedback without
 315 requiring any human annotation. RL allows training on all AUTOPLAY tasks, even those the executor
 316 is not able to solve. However, for better training stability, we restrict RL training to the 8k tasks the
 317 executor can successfully solve at least once. We use GPRO to train the model with RL with group
 318 size 8 across 32 H100 GPUs. See [Appendix C.2](#) for full details.

4 EXPERIMENTS

319 In this section, we demonstrate that AUTOPLAY enables scaling generation of synthetic tasks for
 320 UI agents, that are feasible, diverse and grounded in environment state and functionality all without
 321 human annotations. The resulting high-quality tasks enable scaling of supervised finetuning (SFT)
 322 and reinforcement learning (RL) for post-training MLLMs as capable UI agents. We also show

Method	Model Size	Pass@1
Pass@5		
GPT4o + UI-TARS	43.1	55.3
Seed-VL 1.5	62.1	—
UI-TARS 2 230B	73.3	—
UI-TARS 1.5 7B	26.4	36.2
UI-TARS 72B	37.7	53.0
Qwen-VL 2.5 7B	19.5	27.7
Qwen-VL 2.5 72B	35.0	43.5
AUTOPLAY-7B	40.1 (+20.6 Δ)	58.4 (+30.7 Δ)
AUTOPLAY-72B	47.9 (+12.9 Δ)	68.2 (+24.7 Δ)
AUTOPLAY-3B	34.2	52.2
AUTOPLAY-3B RL	39.9 (+5.7 Δ)	53.4 (+1.2 Δ)

(a) Android World

Method	Pass@1	Pass@5
Pass@5		
o3 + GTA1	34.0	—
AguVis 72B	10.3	—
UI-TARS-1.5 7B	24.5	—
UI-TARS-1.5 72B	42.5	—
Qwen-VL 2.5 7B	3.7	4.1
Qwen-VL 2.5 72B	4.4	5.4
AUTOPLAY-7B	11.4 (+7.7 Δ)	12.1 (+8.0 Δ)
AUTOPLAY-72B	14.5 (+10.1 Δ)	16.0 (+10.6 Δ)

(b) OSWorld

Table 1. Evaluation results on UI agent benchmarks. Pass@1 and Pass@5 values for AUTOPLAY models include in parentheses the change relative to the corresponding base model of the same parameter size.

the tasks synthesized using AUTOPLAY achieve a higher task execution success rate and enable training more capable downstream agents compared to methods that do not explicitly explore the environment (Zhou et al., 2025) and ones that use iterative exploration (Xie et al., 2025; Trabucco et al., 2025). Furthermore, we also evaluate importance of *environment exploration* and *task guidelines* for task generation and find that both components are important to build a effective task proposer.

4.1 AUTOPLAY DATA FOR UI AGENTS

To assess how effectively AUTOPLAY generates training tasks for UI agents, we measure the performance of agents trained with these tasks on established downstream benchmarks: AndroidWorld (Rawles et al., 2024) for mobile agents and OSWorld (Xie et al., 2025) for desktop agents. We use the AUTOPLAY data described in Section 3.3 to train a separate agent per benchmark. In both benchmarks, the agent is evaluated using ground truth success verifier that has access to privileged environment information. We report the average success rate, referred to as Pass@1, over 5 independent random trials since most benchmark tasks have stochastic starting states with different goal details and application content. We also report the Pass@5 metric, which is the percentage of tasks where *any* of the 5 independent trials succeeds.

AUTOPLAY improves agentic capabilities of base models: We use the synthetic data to finetune a base Qwen-2.5-VL-XB into an agentic model AUTOPLAY-XB for a range of sizes $X \in \{3, 7, 72\}$. The resulting models consistently outperform their base as shown in Table 1. For example, AUTOPLAY-7B outperforms its base, Qwen2.5-VL-7B, by 20.6% on AndroidWorld and by 7.7% on OSWorld. These gains demonstrate that the tasks generated by AUTOPLAY are highly relevant to the broad spectrum of UI agent tasks evaluated in both benchmarks, despite requiring *no human data* collection. Similarly, the largest model, AUTOPLAY-72B, surpasses the strong Qwen2.5-VL-72B baseline by 12.9% points on AndroidWorld and 10.1% points on OSWorld, illustrating that even highly capable UI agents benefit substantially from AUTOPLAY’s synthetic tasks. Remarkably, AUTOPLAY-3B achieves a 34.2% success rate on AndroidWorld, nearly matching the performance of the much larger Qwen2.5-VL-72B model, which attains 35.0%.

AUTOPLAY achieves competitive results with Proprietary Baselines trained on human data: Table 1 further shows that AUTOPLAY outperforms strong UI agent models such as UI-TARS-1.5 (Qin et al., 2025), which was trained on large-scale human-annotated GUI data, by 10.2% at the 7B scale and 13.7% at the 72B scale on AndroidWorld. On OSWorld, AUTOPLAY-72B surpasses the two-stage training pipeline used in AguVis-72B (Xu et al., 2024), which leverages both grounding data and human-annotated GUI navigation data, by 5.0% in success rate. Although UI-TARS-2 230B and Seed-VL-1.5 (ByteDanceSeedTeam, 2025) achieve higher performance on AndroidWorld, they rely on substantially larger mixtures of expert models. Similarly, UI-TARS outperforms AUTOPLAY on OSWorld, likely due to its use of curated, human-labeled UI data—whereas AUTOPLAY autonomously explores, proposes tasks, and collects data without human supervision.

AUTOPLAY eventually outperforms the Executor: We also find in Table 1 that AUTOPLAY-72B outperforms the GPT-4o + UI-TARS (OpenAI, 2024; Qin et al., 2025) policy used to collect the data from the generated tasks by 4.8% on the average success rate. This demonstrates the strength of the task verifier for automatically filtering successful trajectories, without ground truth environment

Task Generator	TASK EXECUTION	ANDROID WORLD	
	Pass@1 (↑)	Pass@1 (↑)	Pass@5 (↑)
No Exploration	21.3	28.8 ± 1.5	49.2
Iterative Exploration	56.4	21.6 ± 1.7	33.6
AUTOPLAY w/o task guidelines	43.5	26.7 ± 2.7	38.9
AUTOPLAY	46.0	38.2 ± 3.1	58.5

Table 2. Ablations: We compare AUTOPLAY-7B with two baselines that perform lesser exploration, *No Exploration* and *Iterative Exploration*, as well as AUTOPLAY-7B without using *task guidelines*.

information. Furthermore, AUTOPLAY-7B achieves 12.9% higher Pass@5 than the data collection policy, indicating it's ability to solve new tasks, not just more robustly solve the same tasks.

AUTOPLAY is effective for RL training: In Table 1 we also highlight that it is feasible to perform RL training on the AUTOPLAY generated tasks, according to the process described in Section 3.4. We see a 5.7% gain in AndroidWorld. With RL training, the AUTOPLAY-3B model performs similarly to the AUTOPLAY-7B model trained with just SFT (39.9% versus 40.1% success rate).

4.2 ABLATION ANALYSIS

Next, we look at the importance of individual design decisions in AUTOPLAY. We use AndroidWorld and generate 5,000 Android tasks for all models. We then apply the same executor and verifier described in Section 3.3 to obtain successful trajectories, which are used to fine-tune the Qwen2.5-VL-7B model with supervised learning for AUTOPLAY and variations. Full details are in Appendix F.

Exploration Ablations: To quantify the importance of the environment exploration we compare with two alternative task proposal baselines. The first baseline, *No Exploration*, is based on Zhou et al. (2025); Xie et al. (2025) and, unlike AUTOPLAY, performs no exploration of the domain. Instead, it generates tasks solely from static environment context, such as textual descriptions and application starting screenshots. To implement *No Exploration*, we manually write detailed descriptions of each AndroidWorld app's features and task guidelines in the same style as those used for AUTOPLAY. The second baseline, *Iterative Exploration*, follows Pahuja et al. (2025); Xie et al. (2025) and incorporates limited exploration. However, it only attempts to sequentially execute a series of short-horizon subgoals and summarize them as tasks, without performing broad exploration of the domain. As a result, it is more constrained than AUTOPLAY. Our implementation begins at the home page of the target application. An MLLM proposes a short-horizon subgoal based on the initial screenshot, and the same data collection policy used in AUTOPLAY executes it. The MLLM then proposes the next subgoal, and this process repeats up to 7 times. Finally, a summarizer MLLM takes as input the textual descriptions of the executed subgoals, along with their success or failure, and produces a long-horizon goal for the full trajectory—akin to hindsight relabeling.

AUTOPLAY generates more feasible tasks: Out of the 5k tasks generated by each method, Table 2 shows the same executor succeeds in 46.0% of the AUTOPLAY tasks compared to only 21.3% for *No Exploration* and 56.4% for *Iterative Exploration*. This demonstrates that more tasks generated by AUTOPLAY are feasible. A common failure mode of *No Exploration* baseline was hallucination in the proposed task since it only relies on app functionality description. AUTOPLAY, on the other hand, uses exploration context to ground task details in the actual states of the environment.

AUTOPLAY tasks train better agents: Table 2 also shows that agents trained with AUTOPLAY tasks outperform agents trained with tasks from *No Exploration* and *Iterative Exploration*. Agents trained with AUTOPLAY tasks outperform those trained with *No Exploration* tasks by 9.4% on Pass@1 and Pass@5 success rates on AndroidWorld. AUTOPLAY tends to cover a broader range of functionalities in the environment compared to *No Exploration*. While *Iterative Exploration* does interact with the environment to generate tasks, AUTOPLAY tasks train agents that are 16.6% more successful, as shown by Task Execution Pass@1. This is because *Iterative Exploration* synthesizes long horizon trajectories by stitching easier short horizon subgoals to guide exploration. This consequently leads to less diverse and easier tasks. In contrast, through rounds of long-horizon exploration, AUTOPLAY generates diverse tasks that provide broad coverage over app functionalities.

Exploration generates more diverse tasks: We compare the task distributions generated by AUTOPLAY and *No Exploration* in Figure 3. The distribution is computed over manually defined task categories that cover a broad range of possible tasks. For example, tasks in the *Composition* category

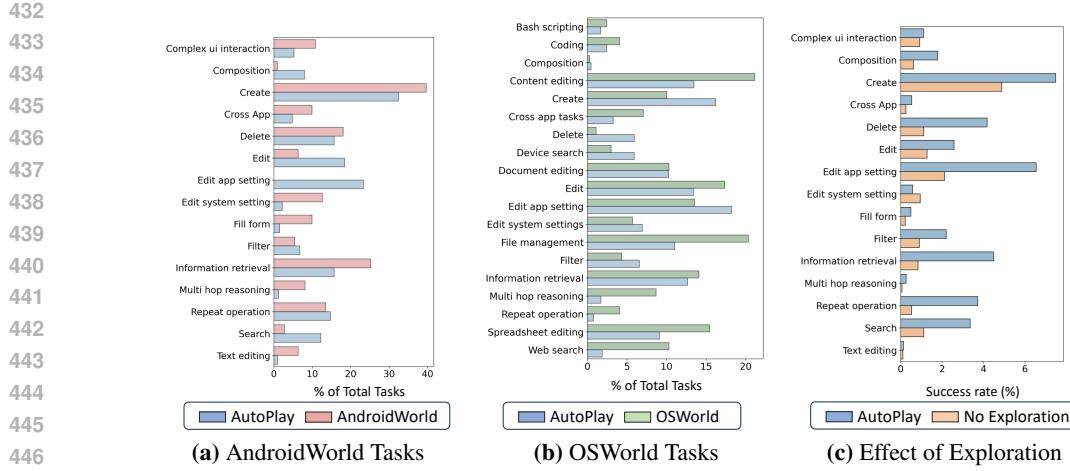


Figure 3. Task Coverage. Left and Middle: For a set of predefined categories, we compare the task distribution across AndroidWorld/OSWorld test tasks and AUTOPLAY generated tasks in light color. **Right:** For a set of predefined categories, we compare the task distribution across AUTOPLAY generated tasks and *No Exploration* generated tasks in light color. Additionally, we show number of tasks that AUTOPLAY-7B executor solves for both task sets in darker color bars.

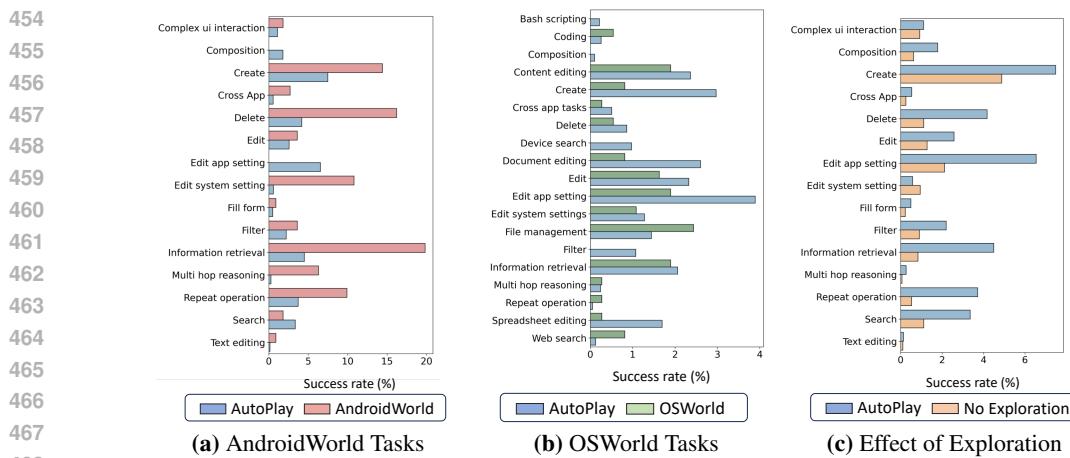


Figure 4. Task Success Rates. Left and Middle: For a set of predefined categories, we compare number of tasks that AUTOPLAY-7B executor solves per task category across AndroidWorld/OSWorld test tasks and AUTOPLAY generated tasks. **Right:** For a set of predefined categories, we compare number of tasks that AUTOPLAY-7B executor solves per task category across AUTOPLAY generated tasks and *No Exploration* generated tasks.

require combining multiple skills or subtasks to achieve the overall goal (e.g., "Find when John is free and schedule a meeting with him for this week."). Additional details about the task categories are provided in Appendix E. Although the two distributions exhibit similar trends—categories that are more prevalent under one method tend to be prevalent under the other—the results consistently show lower execution success rates for *No Exploration* compared to AUTOPLAY, particularly in categories such as deleting, editing, or retrieving in-app data. This highlights the importance of grounding task generation in exploration.

Task guideline ablation: The task guidelines prompts steer the AUTOPLAY task generator towards diverse categories with good domain coverage. Table 2 ablates the impact of the task guideline prompts and illustrates they are important for generating tasks that train performant agents. We see a minor boost in the ability of the executor to solve tasks, showing that task guidelines improve task feasibility. More importantly, guidelines provide a substantial boost in the downstream performance.

App Name	TRAIN ON HOLDOUT APP		HOLDOUT APP	
	Pass@1 (↑)	Pass@5 (↑)	Pass@1 (↑)	Pass@5 (↑)
Markor	26.1	38.4	11.5	38.5
Simple Calendar Pro	33.6	70.5	26.8	47.1
Broccoli-Recipe Expense	20.0	38.6	35.4	46.1
Pro Expense	47.7	66.6	51.1	66.6
Settings	74.4	80.0	72.0	80.0
Average Performance	40.4	58.9	39.4	55.7

Table 3. Generalization Results: Evaluation results of AUTOPLAY-7B to holdout apps *vs.* training on test app.

AUTOPLAY generated tasks cover platform functionality: In addition, we compare of how well AUTOPLAY’s task distribution mimicks a natural task distribution. For this, we use AndroidWorld and OSWorld benchmarks as representative task distributions. As shown in Figure 3, AUTOPLAY produces task covering the majority of the categories. Further, our method seems to follow the distributions of these benchmarks. As shown in Figure 4, for AndroidWorld, the task executor performs quite well on majority of the categories except tasks where skills like fine-grained text editing, complex UI interactions (*e.g.* date-time picker wheels, etc), and cross app navigation is required which leads to the downstream policy perform poorly on these task categories. This clearly suggests, improving the task executor while keeping the task generator the same would lead to better quality synthetic dataset to improve downstream agent performance. For OSWorld, AUTOPLAY covers common task categories like document editing, however, it struggles to generate tasks with cross-app interaction, bash scripting and web search. We attribute this gap in task coverage to lack of sufficiently diverse task guidelines prompts for computer-use domain.

Generalization to unseen environments. In this section, we answer the question: How well do agents trained using AUTOPLAY data generalize to unseen environments compared to training on test environments? To study this question we holdout 5 mobile applications (Settings, Markor, Simple Calendar Pro, Pro Expense, Broccoli-Recipe App) from our training split and evaluate generalization to these apps zero-shot at test time. Due to limited number of environments we train 5 separate models while holding out data from a single environment. In addition, we also train a model only on data collected in the held out environment. We compare the results of zero-shot generalization to held out environment with agents trained only on data from held out environment in Tab. 3. We find that agents trained on larger datasets generalize comparably well to holdout environments compared to training agents on data collected in test environments (39.4% *vs.* 40.4 on Pass@1 success rate).

Verifier analysis. In Tab. 4, we present a human evaluation study to validate effectiveness and human agreement of our verifier. We find the verifiers we used for SFT and RL show \sim 60% accuracy compared to ground truth labels and have high recall. Specifically, we sampled 200 unique tasks with \sim 50% execution success rate and asked 3 human annotators to label them. We label a demonstration as successful if at least 2 annotators agree on it. Next, using human labels as ground truth we compare accuracy, precision and recall for our GPT-4o and Qwen-2.5 VL Instruct model based verifiers in table below. We find both verifiers achieve \sim 60% accuracy and have high recall but low precision. This suggests these models are more optimistic than pessimistic in labelling a task as successful. We find the key failure mode of these models is to mark a long horizon task as successful when an agent completes partial subgoals but doesn’t necessarily complete the full task.

5 CONCLUSION

We introduce AUTOPLAY, a scalable pipeline for synthesizing diverse, feasible, and verifiable tasks for post-training MLLM-based interactive agents. AUTOPLAY achieves this by explicitly exploring interactive environments combined with a carefully curated task guideline prompts to ground task generation in discovered states and accessible functionalities. We demonstrate the effectiveness of AUTOPLAY in training UI agents across both mobile and computer domains. AUTOPLAY generates 20k tasks across 20 Android applications and 10k tasks across 13 applications Ubuntu applications to train mobile-use and computer-use agents. This dataset enables training MLLM-based UI agents that improve in success rates up to 20.0% on mobile-use and 10.9% on computer-use domains. Furthermore, AUTOPLAY generated tasks combined with MLLM reward models enable scaling reinforcement learning training of UI agents, leading to an additional 5.7% gain. Overall, these results establish AUTOPLAY as a scalable approach for post-training MLLM agents, reducing reliance on human annotation while significantly enhancing agent performance across domains.

540 REFERENCES
541

542 Saaket Agashe, Jiuzhou Han, Shuyu Gan, Jiachen Yang, Ang Li, and Xin Eric Wang. Agent S: An
543 Open Agentic Framework that Uses Computers Like a Human. In *International Conference on*
544 *Learning Representations (ICLR)*, 2025. URL <https://arxiv.org/abs/2410.08164>. 1, 3, 30

545 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang,
546 Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan,
547 Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng,
548 Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025. URL
549 <https://arxiv.org/abs/2502.13923>. 6, 15, 25

550 Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai,
551 Lachy Groom, Karol Hausman, Brian Ichter, et al. *pi_0*: A vision-language-action flow model for
552 general robot control. *arXiv preprint arXiv:2410.24164*, 2024. 1

553 ByteDanceSeedTeam. Seed1.5-vl technical report. *arXiv preprint arXiv:2505.07062*, 2025. 7

554 Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and
555 Yu Su. Mind2web: Towards a generalist agent for the web, 2023. 3

556 Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Ayzaan Wahid,
557 Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, et al. Palm-e: An embodied
558 multimodal language model. 2023. 1

559 Linxi Fan et al. Minedojo: Building open-ended embodied agents with internet-scale knowledge. In
560 *arXiv:2206.08853*, 2022. 1

561 Moo Jin Kim, Karl Pertsch, Siddharth Karamchetti, Ted Xiao, Ashwin Balakrishna, Suraj Nair,
562 Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source
563 vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024. 1

564 Jing Yu Koh, Stephen McAleer, Daniel Fried, and Ruslan Salakhutdinov. Tree search for language
565 model agents. *arXiv preprint arXiv:2407.01476*, 2024. 3

566 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
567 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
568 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating
569 Systems Principles*, 2023. 15

570 Wei Li, William Bishop, Alice Li, Chris Rawles, Folawiyo Campbell-Ajala, Divya Tyamagundlu,
571 and Oriana Riva. On the effects of data scale on computer control agents. *arXiv preprint
572 arXiv:2406.03679*, 2024. 1

573 Shikhar Murty, Christopher Manning, Peter Shaw, Mandar Joshi, and Kenton Lee. Bagel: Bootstrap-
574 ping agents by guiding exploration with language, 2024. URL <https://arxiv.org/abs/2403.08140>.
575 3

576 Shikhar Murty, Hao Zhu, Dzmitry Bahdanau, and Christopher Manning. Nnetnav: Unsupervised
577 learning of browser agents through environment interaction in the wild. 2025. 3, 4

578 OpenAI. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. “GPT-4o is an autoregressive
579 omni model ...”. 7

580 OpenAI. Computer-using agent: Introducing a universal interface for ai to interact with the digital
581 world. 2025. URL <https://openai.com/index/computer-using-agent>. 1

582 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni
583 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor
584 Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian,
585 Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny
586 Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks,
587 Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea
588 Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen,

594 Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung,
 595 Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch,
 596 Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty
 597 Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte,
 598 Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh,
 599 Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross,
 600 Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton,
 601 Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton,
 602 Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela
 603 Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan,
 604 Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan
 605 Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt
 606 Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic,
 607 Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung,
 608 Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa
 609 Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov,
 610 Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer
 611 McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob
 612 Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa,
 613 Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mely, Ashvin Nair, Reiichiro Nakano,
 614 Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe,
 615 Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel
 616 Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila
 617 Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle
 618 Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri,
 619 Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl
 620 Ross, Bob Rotstet, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar,
 621 Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard,
 622 Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie
 623 Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie
 624 Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak,
 625 Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick
 626 Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea
 627 Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,
 628 CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave
 629 Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu,
 630 Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan
 631 Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William
 632 Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
 633 6, 14, 20, 25, 26

634 Vardaan Pahuja, Yadong Lu, Corby Rosset, Boyu Gou, Arindam Mitra, Spencer Whitehead, Yu Su,
 635 and Ahmed Hassan Awadallah. Explorer: Scaling exploration-driven web trajectory synthesis for
 636 multimodal web agents. In *Findings of the Association for Computational Linguistics: ACL 2025*,
 637 pp. 6300–6323, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN
 979-8-89176-256-5. URL <https://aclanthology.org/2025.findings-acl.326/>. 3, 8

638 Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John
 639 Wiley & Sons, 2014. 4

640 Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao
 641 Li, Yunxin Li, Shijue Huang, et al. Ui-tars: Pioneering automated gui interaction with native
 642 agents. *arXiv preprint arXiv:2501.12326*, 2025. 1, 6, 7, 20

643 Christopher Rawles, Sarah Clinckemaillie, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Marybeth
 644 Fair, Alice Li, William Bishop, Wei Li, Folawiyo Campbell-Ajala, Daniel Toyama, Robert Berry,
 645 Divya Tyamagundlu, Timothy Lillicrap, and Oriana Riva. Androidworld: A dynamic benchmarking
 646 environment for autonomous agents, 2024. 2, 3, 6, 7, 30

648 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y.K. Li, Y. Wu,
 649 and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language
 650 models, 2024. URL <https://arxiv.org/abs/2402.03300>. 15

651

652 Brandon Trabucco, Gunnar Sigurdsson, Robinson Piramuthu, and Ruslan Salakhutdinov. Insta:
 653 Towards internet-scale training for agents, 2025. 2, 3, 7

654

655 Guanzhi Wang et al. Voyager: An open-ended embodied agent with large language models. In
 656 *arXiv:2305.16291*, 2023. 1

657

658 Xinyuan Wang, Bowen Wang, Dunjie Lu, Junlin Yang, Tianbao Xie, Junli Wang, Jiaqi Deng, Xiaole
 659 Guo, Yiheng Xu, Chen Henry Wu, Zhennan Shen, Zhuokai Li, Ryan Li, Xiaochuan Li, Junda Chen,
 660 Boyuan Zheng, Peihang Li, Fangyu Lei, Ruisheng Cao, Yeqiao Fu, Dongchan Shin, Martin Shin,
 661 Jiarui Hu, Yuyan Wang, Jixuan Chen, Yuxiao Ye, Danyang Zhang, Dikang Du, Hao Hu, Huarong
 662 Chen, Zaida Zhou, Haotian Yao, Ziwei Chen, Qizheng Gu, Yipu Wang, Heng Wang, Diyi Yang,
 663 Victor Zhong, Flood Sung, Y. Charles, Zhilin Yang, and Tao Yu. Opencua: Open foundations for
 664 computer-use agents, 2025. URL <https://arxiv.org/abs/2508.09123>. 1

665

666 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le,
 667 and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.
 668 URL <https://arxiv.org/abs/2201.11903>. 6, 25

669

670 Jingxu Xie, Dylan Xu, Xuandong Zhao, and Dawn Song. Agentsynth: Scalable task generation for
 671 generalist computer-use agents. *arXiv preprint arXiv:2506.14205*, 2025. 2, 3, 7, 8, 28

672

673 Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing
 674 Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio
 675 Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. Osworld: Benchmarking multimodal agents
 676 for open-ended tasks in real computer environments, 2024. 2, 3, 6, 30

677

678 Yiheng Xu, Zekun Wang, Junli Wang, Dunjie Lu, Tianbao Xie, Amrita Saha, Doyen Sahoo, Tao Yu,
 679 and Caiming Xiong. Aguvis: Unified pure vision agents for autonomous gui interaction, 2024. 7

680

681 Yan Yang, Dongxu Li, Yutong Dai, Yuhao Yang, Ziyang Luo, Zirui Zhao, Zhiyuan Hu, Junzhe Huang,
 682 Amrita Saha, Zeyuan Chen, Ran Xu, Liyuan Pan, Caiming Xiong, and Junnan Li. Gta1: Gui
 683 test-time scaling agent, 2025. URL <https://arxiv.org/abs/2507.05791>. 5, 6, 20, 30

684

685 Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable
 686 real-world web interaction with grounded language agents. In S. Koyejo, S. Mohamed, A. Agarwal,
 687 D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, vol-
 688 umne 35, pp. 20744–20757. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/82ad13ec01f9fe44c01cb91814fd7b8c-Paper-Conference.pdf. 1

689

690

691 Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. τ -bench: A benchmark for
 692 tool-agent-user interaction in real-world domains, 2024. URL <https://arxiv.org/abs/2406.12045>.
 693 1, 3

694

695 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,
 696 Yonatan Bisk, Daniel Fried, Uri Alon, et al. Webarena: A realistic web environment for building
 697 autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023. URL <https://webarena.dev>. 1, 3

698

699 Yifei Zhou, Qianlan Yang, Kaixiang Lin, Min Bai, Xiong Zhou, Yu-Xiong Wang, Sergey Levione,
 700 and Erran Li. Proposer-agent-evaluator (PAE): Autonomous skill discovery for foundation model
 701 internet agents. In *ICML*, 2025. URL <https://arxiv.org/abs/2412.13194>. 2, 3, 7, 8

702	Verifier Model	Accuracy (\uparrow)	Precision (\uparrow)	Recall (\uparrow)
703	GPT-4o	59.3	56.9	82.8
704	QwenVL-2.5 Instruct	61.8	57.7	95.9

706 **Table 4. Verifier analysis:** Human evaluation of verifier accuracy on AndroidWorld data.
707

708	Verifier	Pass@1 (\uparrow)	Pass@5 (\uparrow)	Dataset Size
709	1) GPT-4o	40.1	58.4	8k
710	2) Qwen2.5-VL 72B	38.6	54.5	12k

712 **Table 5. Verifier Ablations.** Evaluation results of training Qwen-VL 2.5 7B using SFT with different verifiers
713 on AndroidWorld.
714715

A USE OF LLM FOR WRITING

718 We use LLMs to assist with specific aspects of paper writing which includes: fix grammatical errors,
719 sentence polishing, and paraphrasing parts of for sections 1, 3 and 4. LLMs were used for: (1)
720 grammar checking and improving clarity and readability of text, (2) suggestions for writing certain
721 sentences via paraphrasing. For each section the original content was written by the authors and
722 edited using GPT-4o [OpenAI et al. \(2024\)](#) LLM. Followed by this round of LLM-assisted editing, the
723 authors verified and polished the content written by the LLM to fix issues with generated text and
724 create the final text used for the paper.

725

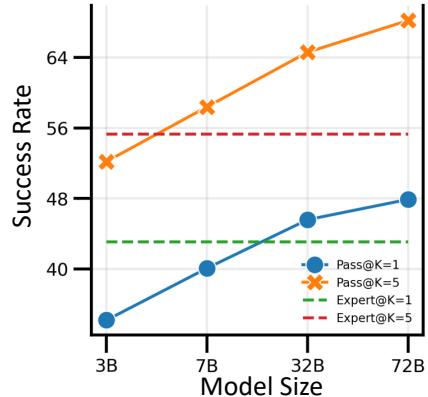
B ADDITIONAL EXPERIMENTS

728 **Impact on distillation performance with different ver-729 ifiers.** In this section, we ablate the choice of verifier
730 model used for data filtering for SFT and its impact on
731 downstream performance. Specifically, we compare use of
732 proprietary MLLMs like GPT-4o vs. open-source MLLM
733 Qwen2.5-VL Instruct. This experiment is especially im-
734 portant to evaluate what capability level is required in
735 a MLLM to be used as a effective verifier and whether
736 smaller open-source MLLMs can be used as effective ver-
737 ifiers to enable RL scaling as proprietary models cannot be
738 used due to cost implications. To do so, we use all $\sim 20k$
739 tasks and demonstrations generated by our synthetic task
740 proposer and executor and verify these using both GPT-4o
741 and Qwen2.5-VL Instruct 72B to create two SFT datasets.
742 In Tab. 5, we present results of finetuning Qwen2.5-VL
743 Instruct 7B model on these trajectory datasets and present
744 evaluation results on AndroidWorld benchmark. We find
745 that model trained on GPT-4o verified trajectories out-
746 performs Qwen2.5-VL Instruct 72B verified trajectories,
747 however, the difference in performance is quite small 1.5% suggesting smaller open-source models
748 are capable verifiers that can enable RL training in mobile-use domains.

748

B.1 GENERALIZATION RESULTS

750 So far for all experiments we have used same set of training and test environments. Even though these
751 experiments demonstrate effectiveness of our pipeline in providing a broad coverage over training
752 task distribution and its impact on downstream performance, another important question is: How well
753 do agents trained using synthetic data generalize to unseen environments compared to training on test
754 environments? To study this question we holdout 5 mobile applications (Settings, Markor, Simple
755 Calendar Pro, Pro Expense, Broccoli-Recipe App) from our training split and evaluate generalization
to these apps zero-shot at test time. Due to limited number of environments we train 5 separate

756 **Figure 5.** Effect of increasing model size in
757 AUTOPLAY.

756 models while holding out data from a single environment. In addition, we also train a model only
 757 on data collected in the held out environment. We compare the results of zero-shot generalization
 758 to held out environment with agents trained only on data from held out environment in ???. We find
 759 that agents trained on larger datasets generalize comparably well to unseen (or test) environments
 760 compared to training agents on data collected in test environments.
 761

762 C TRAINING DETAILS

764 C.1 ALGORITHM PARAMETERS

766 We provide parameter values for AUTOPLAY, as listed in Algorithm 1, for both AndroidWorld and
 767 OSWorld experiments:

769 Parameter Name	770 AndroidWorld	771 OSWorld
770 S # of apps	771 20	772 13
771 K # of exploration turns	772 3	773 5
772 P # of task guidelines	773 4	774 4
773 K # of task generations per guideline and context	774 50	775 50

776 C.2 RL TRAINING DETAILS

777 We use GRPO (Shao et al., 2024) for our RL finetuning experiments using 32 H100 GPUs. As
 778 we require a MLLM verifier for evaluating task success and reward generation we use 2 GPUs
 779 on everywhere GPU node to host the Qwen2.5-VL Instruct 32B (Bai et al., 2025) MLLM using
 780 vLLM (Kwon et al., 2023) inference engine. The verifier only takes the final 8 frames of task
 781 execution as input and predicts a binary evaluation of the task execution in addition to reasoning and
 782 task summary as described in Sec. 3.3.

783 Each GPU node runs environment workers and RL trainer on 6 of the GPUs with 8 environments per
 784 GPU. We evaluate the policy after 120 GRPO updates. This is equivalent to $8 \times 24 \times 16 \times 120 =$
 785 368,640 environment samples. Additional hyperparameters are detailed in Tab. 6.

787 D ADDITIONAL AUTOPLAY DETAILS

789 D.1 TASK PROPOSER AGENT

791 AUTOPLAY task proposer agent operates in two stages where it first explores the environment in
 792 a goal-agnostic manner to collect experience from the environment that can aid in synthesis of
 793 high-quality task that are grounded in environment state and feasible to execute. The task proposer
 794 takes as input the task proposer prompt, goal-agnostic exploration observations and task guideline
 795 prompt to generate a sequence of tasks. We present the task proposer prompt prefix in Appendix D.1.

796 Next, for mobile-use domain we define 4 separate task guideline prompts that are specific to mobile-
 797 use domain which the task proposer uses in addition to exploration experience to generate tasks.
 798 These prompts cover tasks of following types: (1.) Feature-Use Appendix D.1: Tasks that use basic
 799 features of the app and provide broad coverage over all features shown, (2.) Feature-Composition Ap-
 800 pendix D.1: Tasks that composes multiple feature-use tasks to create complex tasks with multiple
 801 subtasks, (3.) Information Retrieval Appendix D.1: Tasks that require searching for specific infor-
 802 mation requiring or answering queries asked by users about state of the environment, (4.) Feature-
 803 Repetition Appendix D.1: Tasks that require executing feature-use tasks over multiple entities stored
 804 in an application (e.g. deleting multiple calendar events).

805 For computer-use domain we define 2 separate task guideline prompts mentioned in Appendix D.1
 806 and Appendix D.1.

808 **Task Generator MLLM Prompt (Mobile-Use: Android and Computer-Use: Ubuntu)**

Parameter	Value
Number of GPUs	24
Number of environments per GPU	8
Rollout length	16
GRPO Group Size	8
Number of mini-batches per epoch	4
LR	$1.e^{-6}$
Entropy coefficient	0.0
KL divergence coefficient	0.0
LR scheduled	10 linear warmup steps from 0 LR to 1e-6 LR
Max gradient norm: 1.0	
Optimizer	Adam
Weight decay	0.0

Table 6. Hyperparameters used for RL finetuning using GRPO.

You are a capable UI understanding agent. You will be provided a set of images from the {PLATFORM} app that shows different features and the current state of the app. Your task is to convert the described functionalities into a list of tasks that a useful UI understanding agent should be able to complete. Use the images from the app to ground these tasks to the ones that are feasible. Propose as many diverse tasks as possible to cover broad range of features.

For all described and demonstrated functionalities output a list of up to {NUM_TASKS} unique tasks that can be executed in the app as a JSON which should be in the following format:

```

tasks = [
    {
        "thought": "<Detailed thought and reasoning for the proposed task, why is the task simple enough to execute and whether it satisfies the task proposal guidelines>",
        "instruction": <natural language instruction describing a task with/without template
params with name of the app>,
        "tag": <few words describing the type of task>,
        "app_name": <name of the app>,
        "template_params": {
            "param_name": {
                "description": <param description>,
                "possible_values": [<list of 5 random values>]
            }
        }
    }
]
{TASK_GUIDELINE_PROMPT}
{ENVIRONMENT_CONTEXT}

```

Feature-Use Task Guidelines Prompt (Mobile-Use: Android)

For each functionality/feature make sure to follow these guidelines while generating tasks:

1. Create all possible tasks. For example, if there is a clock app on the phone which has worldwide clock feature then you can create tasks like: "Add worldwide clock for <city>", "Remove worldwide clock for <city>", "Add worldwide clock for <city_1, city_2> and reorder to put <city_2> before <city_1>". Make sure to add the <param> details in "template params" in JSON to allow creation of diverse tasks.
2. In any instruction if there are going to be templated parameters then add it in following format: {param_name}
3. If a task requires creating/adding/editing/deleting any information/event/entry make sure all required entities for completing the task are parameterized/templated. For example:
 - a.) For a task that requires creating notes both the name of the new note and content of the note should be specified as parameters.
 - b.) Similarly, if a task requires editing some details describe the exact edit you'd like to make, what entity needs to be edited. When proposing such tasks, only ask edits to existing entities shown in the images from the app for such tasks. If no such entities exist ask the agent to first create one and then edit it.
 - c.) If a task requires deleting an entity, describe the entity that needs to be deleted. When proposing such a task, only ask to delete an entity that already exists as shown in the images from the app for such tasks. If no such entities exist ask the agent to first create one and then delete it.

864 d.) If a task requires copying an entity, describe the entity that needs to be copied and
865 where it should be copied to (ex: which folder or date). When proposing such a task, only ask
866 to copy an entity that already exists as shown in the images from the app for such tasks. If
867 no such entities exist ask the agent to first create one and then copy it.
868 e.) For all parameters that are templated if any parameter refers to an entity for edit,
869 delete or copy task make sure that the list of possible values only contain entities that
870 exists or will be created.
871 4. Ensure that instructions are unambiguous and clearly describes a actual task that can be
872 performed on the app.
873 5. Do not include tasks that require accessing/uploading/capturing content from real world. For
874 example, scanning a document, recording a new video, taking a picture using camera.
875 6. Task instructions should be natural user requests someone might actually ask a capable UI agent
876 to complete on the app.
877 7. If any supports browsing information from the web or library make sure to include that in the
878 task instruction. For example, if a task requires searching for a book on the web, make sure to
879 include that in the task instruction.

Information Retrieval Task Guidelines Prompt (Mobile-Use: Android)

For each proposed task make sure to follow these guidelines:

1. Create tasks that require retrieving different types of information that are useful for a user to make decisions. These tasks should be natural user requests someone might ask in real world. For example, if there is a calendar app on the phone you can propose tasks like: "How many meetings do I have scheduled for {date}?", "What meetings do I have to attend between {start_time} and {end_time} on date {date}? List each event separated by comma". Make sure to add each templated param {param} in "template params" in JSON to allow creation of diverse tasks.
2. Task instructions should be natural user requests someone might actually ask a capable UI agent to complete on the app. For example, "Check the current sample rate set in the Audio Recorder app." is a bad task because of how it is specified. Instead, propose tasks like "What is the current sample rate set in the Audio Recorder app?" or "Can you tell me the current sample rate set in the Audio Recorder app?". The latter two tasks are better because they explicitly ask for the information to be retrieved and are more natural user requests.
3. For all such tasks, also generate a "answer" field in the task JSON which contains a natural language answer to the task. This answer should be a valid answer that can be retrieved from the app. For example, if the task is "What is the weather forecast for {city} on {date}?", the answer should be a valid weather forecast for the specified city and date.
4. For any task instruction if there are going to be templated parameters then add it in following format: {param_name}
5. For all such task propose different varieties of information retrieval tasks that require searching for information in the app and covers different features. Here are some examples for different apps:
 - a.) For a notes app, you can ask "How many todos I have listed in the notes app?", "What todos are pending for today in the list from note app?",etc.
 - b.) For a fitness app, "What workouts do I have planned for this week?", "How long did I workout in last week?", etc.
6. Only proposes tasks where the last action the agent should require to solve the task is to return a natural language response with the information requested to the user. For example, "What is the weather forecast for {city} on {date}?" is a good task because it requires the agent to return a natural language response with the weather forecast for the specified city and date. Whereas, "How do I view the weather forecast?" or "How can I set the city to {city} in weather app?" is not a good task because it does not specify the information to be retrieved and requires the agent to show a demonstration in the app to view the weather forecast.
7. Do not propose ambiguous tasks that do not specify the information to be retrieved. For example, "What is the weather like?" is too generic and should be avoided. Instead, propose tasks like "What is the weather forecast for {city} on {date}?".
8. If any app supports browsing information from the web or library make sure to include tasks that require searching information from web. For example, "How much is the price of a book titled {book}?", etc.

Feature Composition Task Guidelines Prompt (Mobile-Use: Android)

908 For each task proposed make sure to follow these guidelines:
909 1. Each task should be a composition of multiple subtasks that require a UI agent to execute
910 sequence of subtasks across multiple functionalities/subtasks. For example, if there is a clock
911 app on the phone which has worldwide clock, alarms, and timer feature then you can create tasks
912 like: "Add worldwide clock for {city} then set an alarm for {time} on days {}", "Remove worldwide
913 clock for {city} and set a timer for {hour} hours, {minute} minutes", "Delete the alarm for {time}
914 and delete all the world clocks". Make sure to add each templated {param} details in "template
915 params" in JSON to allow creation of diverse variants of each tasks.
916 2. Task instructions should be natural user requests, unambiguous and clearly describes a actual
917 task that can be performed on the app. For example, "Create a calendar event for meeeting titled
918 {title} at {time} on {date} for duration {duration} for meeting with Bob then delete the first
919 event on {date2}" is a good task as it specifies all required details. In contrast, a task like
920 "Create a calendar event for titled {title} on {date} then delete the first event on {date2}" is a
921 bad task as it does not specify the start time/duration of the event.

918 2. For any instruction if there are going to be templated parameters then specify it in following
 919 format in the instruction: {param_name}
 920 3. For all such tasks that require creating/adding/editing/deleting any information/event/entry
 921 make sure all required entities for completing the task are parameterized/templated. For example:
 922 a.) For notes app, "Create a note titled {note_title} with content {note_content} in {folder}
 923 folder and then delete note titled {note_title_2}" is a good task as it templates the title,
 924 content, and location of notes.
 925 b.) For expense app, "Create expenses for {expense_name_1}, amount {amount_1}, category
 926 {category_1}, note {note_1} followed by expense {expense_name_2}, amount {amount_2}, category
 927 {category_2}, note {note_2} and then delete all duplicate expenses".
 928 c.) For files app, you can ask "Search for file named {name_1} and then go delete the files
 929 named {name_2}, {name_3}".
 930 4. Do not include tasks that require accessing/uploading/capturing content from real world. For
 931 example, scanning a document, recording a new video, taking a picture using camera.
 932 5. If any app supports browsing information from the web or library make sure to include that in
 933 the task instruction. For example, if a task requires searching for a book on the web, make sure
 934 to include that in the task instruction.

931 Subtask Repetition Task Guidelines Prompt (Mobile-Use: Android)

932 For each task proposed make sure to follow these guidelines:
 933 1. Each task instruction you propose should repeatedly ask to execute the same feature or subtask
 934 for a single functionality. For example, if there is a calendar app on the phone which has
 935 multiple calendar events then you should propose tasks like: "Delete events {event_1}, {event_2},
 936 {event_3}", "Delete all events on date {date_1} {date_2}", "Delete events {event_1} on {date_1},
 937 {event_2} on {date_2}, {event_3} on {date_3}" and so on. Make sure to add each templated {param}
 938 details in "template params" in JSON to allow creation of diverse variants of each tasks. Make
 939 sure the for each instance of event in template params unique values from screenshots shown are
 940 used.
 941 2. Task instructions should be natural user requests, unambiguous and clearly describes a actual
 942 task that can be performed on the app. The instructions can specify same task in different ways.
 943 For example, "Delete all events on {date_1}" can also be specified as "Delete all events on {day}
 944 of the {week}", and a task like "Delete all events on this weekend" can be specified as "Delete
 945 all events on Saturday and Sunday of current week".
 946 3. Also specify tasks in various ways that require inferring details of the entity being referred
 947 by looking at details from the screenshot. For example, if the app is a expense app you can
 948 specify tasks like "Delete expenses that have expense amount greater than {amount}" or "Delete all
 949 expenses that are in {category} category" or "Delete all expenses that have a note containing
 950 {note_content}".
 951 4. Propose tasks that require repeatedly executing same feature for all types of tasks like
 952 create/delete/edit. For example, if the app is a notes app you can propose tasks like "Create
 953 notes titled {note_1}, {note_2}, {note_3} with content {content_1}, {content_2}, {content_3}",
 954 "Delete notes titled {note_1}, {note_2}, {note_3}", "Edit notes titled {note_1} to change content
 955 to {new_content_1}, edit note titled {note_2} to change content to {new_content_2}" and so on.
 956 5. For any instruction if there are going to be templated parameters then specify it in following
 957 format in the instruction: {param_name}
 958 6. Do not include tasks that require accessing/uploading/capturing content from real world. For
 959 example, scanning a document, recording a new video, taking a picture using camera.
 960 7. If any app supports browsing information from the web or library make sure to include that in
 961 the task instruction. For example, if a task requires searching for a book on the web, make sure
 962 to include that in the task instruction.

956 Feature-Use Task Guidelines Prompt (Computer-Use: Ubuntu)

957 For each functionality/feature make sure to follow these guidelines while generating tasks:
 958 1. Create all possible tasks that use features shown by the demonstrations. For example, if there
 959 is a code IDE app on the desktop then you can propose tasks like: "Create a new project at path
 960 <path> from UI or terminal", "Install an extension <extension_name> for auto formatting text",
 961 "Change the settings of my editor to set line-length to <value>", etc. Make sure to add the
 962 <param> details in "template params" in JSON to allow creation of diverse tasks.
 963 2. In any instruction if there are going to be templated parameters then add it in following
 964 format: {param_name}
 965 3. If a task requires creating/adding/editing/deleting any information/event/entry make sure all
 966 required entities for completing the task are parameterized/templated. For example:
 967 a.) For a task that requires creating a new file in a coding IDE specify the name of the file,
 968 location, and content of the file as parameters. Similarly, if a task requires creating a new
 969 project specify the name of the project, location, and any other required details as
 970 parameters.
 971 b.) If a task requires editing some details describe the exact edit you'd like to make, what
 972 entity needs to be edited. When proposing such tasks, only ask edits to existing entities
 973 shown in the images from the app for such tasks. If no such entities exist ask the agent to
 974 first create one and then edit it. For example, if a task requires editing a image in image
 975 editing app, specify the image to be edited, the edit to be made, and any other required
 976 details as parameters.

972 c.) If a task requires deleting an entity, describe the entity that needs to be deleted. When
 973 proposing such a task, only ask to delete an entity that already exists as shown in the images
 974 from the app for such tasks. If no such entities exist ask the agent to first create one and
 975 then delete it.
 976 d.) If a task requires copying an entity, describe the entity that needs to be copied and
 977 where it should be copied to (ex: which folder or date). When proposing such a task, only ask
 978 to copy an entity that already exists as shown in the images from the app for such tasks. If
 979 no such entities exist ask the agent to first create one and then copy it.
 980 e.) For all parameters that are templated if any parameter refers to an entity for edit,
 981 delete or copy task make sure that the list of possible values only contain entities that
 982 exists or will be created.
 983 4. Ensure that instructions are unambiguous and clearly describes a actual task that can be
 984 performed on the app.
 985 5. Task instructions should be natural user requests someone might actually ask a capable UI agent
 986 to complete on the app.
 987 6. If any supports browsing information from the web or library make sure to include that in the
 988 task instruction. For example, if a task requires searching for a book on the web, make sure to
 989 include that in the task instruction.

990 **Feature Composition Task Guidelines Prompt (Computer-Use: Ubuntu)**

991
 992 Available Primitives:
 993 - search: searches for something using the search bar. You can optionally filter the results based
 994 on the given criteria.
 995 - filter: filters the results based on the given criteria. You can compose multiple filters to
 996 form a single filter. A filter can be direct, or it needs to be inferred using multi_hop_reasoning.
 997 - edit: edit/modifies a property (changing to celsius, sorting results, comparing, modifying the
 998 view by clicking on a button, etc.)
 999 - delete: deletes something from the page. can also be used to delete something from the cart.
 1000 - add: add something to the cart. You can combine multiple add primitives, you can also add
 1001 multiple quantities of the same item.
 1002 - multi_hop_reasoning: criteria for selection/filtering is not mentioned directly, but requires
 1003 reasoning to be performed. Examples: "one month from now", "30% of an income of \\$8000".
 1004 - repeat: repeats the results based on the given criteria. Example: add red wine, white wine. Add
 1005 2 bottles.
 1006 - navigate: navigates to the given URL, click on a link to navigate to its page. select one of the
 1007 entries from a list.
 1008 - form: fills out a form (contact, login, enter numbers to perform a calculation, filling in zip
 1009 code, etc.). We can stack multiple forms to create a task.
 1010 - and: logical operation
 1011
 1012 For each functionality/feature make sure to follow these guidelines while generating tasks:
 1013 1. Generate 10 unique tasks each with 2, 4, 6 primitives in your task composition (minimum 3,
 1014 maximum 6). Make sure the tasks are diverse and use diverse composition of primitives. Always
 1015 include atleast a few tasks which uses multi_hop_reasoning, repeat, and add primitives.
 1016 2. Task instructions should be natural user requests someone might actually ask a capable UI agent
 1017 to complete on the app. Tasks should be motivated by daily use cases.
 1018 3. For any instruction if there are going to be templated parameters then add it in following
 1019 format: {param_name}
 1020 4. Ensure that instructions are unambiguous and clearly describes an actual task that can be
 1021 performed.
 1022 5. Create realistic, executable tasks that combine multiple primitives logically in a meaningful
 1023 sequence.
 1024 6. Include primitive composition as a function-like string (e.g., "search(product,
 1025 filter=[filter(price, 50), filter(location, NYC)])")
 1026 7. For all tasks that require creating/adding/editing/deleting any information/entity make sure
 1027 all required entities for completing the task are parameterized/templated. For example:
 1028 a.) For tasks that require adding items, both the item details and quantities should be
 1029 specified as parameters.
 1030 b.) For tasks that require filtering, the filter criteria should be templated to allow diverse
 1031 task variants.
 1032 c.) For tasks requiring multi-hop reasoning, the reasoning criteria should be clearly
 1033 parameterized.
 1034 8. IMPORTANT: Do not include tasks that require accessing/uploading/capturing content from real
 1035 world. For example, scanning documents, recording videos, taking pictures using camera.
 1036 9. IMPORTANT: Do not generate tasks which require login.
 1037 10. Tasks should demonstrate composition of primitives working together, not just sequential
 1038 execution of unrelated actions.
 1039 11. Include template parameters where appropriate for task variation, with at least 5 possible
 1040 values for each parameter.
 1041 12. Use the provided image context to ground the task in the screenshots. Don't hallucinate any
 1042 template parameters.

1026
1027

D.2 EXPLORATION AGENT

1028 Our explorer agent is instantiated using the implementation of the task executor agent mentioned
 1029 in Appendix D.3 with two key differences. First, the task instruction for each exploration turn is
 1030 a generic exploration goal of the format “Explore the APP_NAME app exhaustively to access all
 1031 features, functionalities and data stored on the app.”. In addition, across multiple runs of exploration
 1032 the explorer agent maintains an explicit memory of environment state from the past exploration
 1033 turns in the same application. This memory is a text representation and it is only kept in context
 1034 when running multiple exploration turns for the same application. In order to convert a exploration
 1035 trajectory to the memory representation we use GPT-4o [OpenAI et al. \(2024\)](#) as our summarizer
 1036 MLLM. Specifically, we give the sequence of observations from the exploration trajectory as input to
 1037 the summarizer MLLM and ask it to output a structured response that describes the functionalities of
 1038 the environment observed by they agent in the exploration turn followed by a description of data or
 1039 configs stored in the application that would be relevant for task curation. The prompt used for the
 1039 summarizer MLLM is described in Appendix D.2.

1040
1041**Summarizer MLLM Prompt (Mobile-Use: Android and Computer-Use: Ubuntu)**

1042 You are a capable UI understanding agent. You will be provided a sequence of images from an app or
 1043 a website that shows how to access different features and the data currently stored in the app.
 1044 Your task is to summarize the actions taken and features, functionalities navigated and interacted
 1044 with by the user in detail.

1045 Your task is to output the information in following format:
 1046 {
 1047 "action summary": "A bulleted list of actions taken and pages visited to explore the features
 1048 and functionalities explored by the agent that can enable generation of meaningful UI control
 1049 tasks. Describe the pages and features visited in the order they were explored."
 1050 "data stored": "A bulleted list of data found on the explored by the agent that can enable
 1051 generation of meaningful UI control tasks. Describe the data uncovered in the order of the
 1052 features visited as they were explored."
 1053 }

1054

D.3 TASK EXECUTOR AGENT

1055

1056 Our task executor agent is instantiated as a modular agent that decomposes the task execution into
 1057 high-level task planning using a high-level planner MLLM and low-level action execution using
 1058 grounding model. In addition, the high-level planner uses a reflection tool which takes as input the
 1059 past observation, past action taken, and current observation to describe the transition caused by the
 1060 action. The high-level planner MLLM takes as input the task instruction, current observation, history
 1061 of past actions, reflection trace of past action, and outputs a high-level action. If the high-level action
 1062 are click-based interaction then it uses the grounding model to localize the coordinate location.

1063

1064 For mobile domains, we use GPT-4o ([OpenAI et al., 2024](#)) as both the planner and reflection model,
 1065 and UI-TARS-1.5 7B ([Qin et al., 2025](#)) as the grounding model. We provide the prompt used by the
 1066 high-level planner in Appendix D.3, the reflection tool prompt in Appendix D.3, and the prompt for
 1067 UI-TARS-1.5 7B grounding model in Appendix D.3.

1068

1069 For desktop domains, we use GPT-4o as the planner and reflection model, and GTA1-7B ([Yang](#)
 1070 [et al., 2025](#)) as the grounding model. We provide the prompt used by the high-level planner in Ap-
 1071 pendix D.3, the reflection tool prompt in Appendix D.3, and the prompt for GTA1-7B grounding
 1072 model in Appendix D.3. In the desktop setting, we additionally supply the high-level policy with
 1073 heuristic actions (detailed in Appendix G) to improve data collection policy success rate.

1073

1074

1075
1076**High-Level Planner MLLM Prompt (Mobile-Use: Android)**

1077 You are a capable GUI assistant designed to help users navigate and interact with mobile
 1078 applications. At the beginning of each task, you will be provided with a natural language
 1079 description of the task.

1077

1078

1079

Then, at each step, you will be provided with:

1. the screen image.
2. the history of actions taken on the environment and any feedbacks from an expert critic

```

1080 3. feedback on the last action taken.
1081
1082 Your task is to analyze the goal, current screen, history of actions and feedback about the last
1083 action, think step-by-step and generate a natural language plan that clearly describes the next
1084 action with relevant details element or area on the screen for interaction. Finally, also output
1085 the next action described in natural language a single sentence. Only plan for next action using
1086 actions described above. Drag action is not supported, use clicks instead of drag.
1087
1088 Here are some more guidelines:
1089 1. If the task has been completed, you should call the terminate action.
1090 2. Avoid getting stuck in trying to call the same action over and over again. If you think you are
1091 stuck, try to find alternative ways to complete the task.
1092 3. Try to accomplish the task with the least number of actions.
1093 4. If the previous action failed, take corrective action by taking an alternative action.
1094
1095 Goal: {TASK_INSTRUCTION}
1096
1097 Action History:
1098 {ACTION_HISTORY}
1099
1100 Critic Response for Last Action:
1101 {REFLECTION_LLM_RESPONSE}
1102
1103 Summary of screen changes:
1104 {TRANSITION_SUMMARY}
1105
1106 Instructions: Based on the goal, current screen, history of actions containing feedback about the
1107 last action, think step-by-step and provide the plan and next action.
1108
1109 At each timestep output the next action in following format:
1110 {
1111   reason: "Step-by-step thinking for the action to be taken.",
1112
1113   action: "The action to be taken. Choose one of the available actions. tap_on_element if you
1114   want to tap on an element, long_press_on_element if you want to long press on a
1115   element/location, type_text_in_element if you want to type text in an element, scroll_screen
1116   if you want to scroll the screen, answer if you want to answer the question asked by the user,
1117   terminate if you want to terminate the session, navigate_back if you want to navigate to
1118   previous screen, navigate_home if you want to navigate to the home screen, wait if you want to
1119   wait for 5 seconds before continuing, open_app if you want to open a specific app.",
1120
1121   element_description: "A descriptions containing visual and textual details as well as spatial
1122   location of the UI element to be interacted with. Only set the element description if the
1123   action is tap_on_element, long_press_on_element, or type_text_in_element, otherwise set the
1124   element description to be an empty string. Examples of descriptions: circular design with a
1125   pattern resembling a camera shutter or a mandala, in white on a light brown background.
1126   Otherwise set the element description to be an empty string. Select the checkbox next to the
1127   label 'Remember me'.",
1128
1129   text: "The text to be entered in the UI element. Only set with value to type in the text
1130   element.",
1131
1132   direction: "The direction to scroll the screen. Can be one of [UP, DOWN, LEFT, RIGHT].",
1133
1134   answer: "The answer to the question asked by the user. Only set this field if the action is
1135   answer, otherwise set the answer to be an empty string.",
1136
1137   app_name: "The name of the app to be opened. Only set this field if the action is open_app,
1138   otherwise set the app_name to be an empty string."
1139 }
1140
1141 Current Observation:
1142 {OBSERVATION}
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3020
3021
3022
3023
3024
3025
3026
3027
3028
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3070
3071
3072
3073
3074
3075
3076
3077
3078
3079
3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185
3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
321
```

```

1134 Your task is to think carefully and analyze the current action, as well as the screen before and
1135 after the action. Use this information to describe the changes in the screen and provide feedback
1136 about whether
1137 the action was successful or not.
1138 If the action changes state on the screen but does not open a new view, focus on the element which
1139 was effected to assess success.
1140 If the action was a scroll or swipe and the UI elements did not change, the action likely failed.
1141 If the action was to type in the text field, the evaluation should be whether or not the textfield
1142 has the exact text typed in.
1143 Don't tell the human what to do, just provide feedback on whether the action was successful or not.
1144
1145 Goal: {TASK_INSTRUCTION}
1146
1147 Observation before action:
1148 {PREVIOUS_OBSERVATION}
1149
1150 Action executed:
1151 {ACTION}
1152
1153 Observation after action:
1154 {CURRENT_OBSERVATION}

```

Low-Level Planner MLLM Prompt - UI-TARS 1.5 7B (Mobile-Use: Android)

```

1151 You are a helpful assistant.
1152
1153 You are a GUI agent. You are given a task and your action history, with screenshots. You need to
1154 perform the next action to complete the task.
1155 ## Output Format
1156 ...
1157 Thought: ...
1158 Action: ...
1159 ...
1160 ## Action Space
1161 click(point='<point>x1 y1</point>')
1162 long_press(point='<point>x1 y1</point>')
1163 type(content='') #If you want to submit your input, use "\n" at the end of `content`.
1164 scroll(point='<point>x1 y1</point>', direction='down or up or right or left')
1165 open_app(app_name='')
1166 drag(start_point='<point>x1 y1</point>', end_point='<point>x2 y2</point>')
1167 press_home()
1168 press_back()
1169 finished(content='xxx') # Use escape characters '\\', '\\', and \\n in content part to ensure we can
1170 parse the content in normal python string format.
1171
1172 ## Note
1173 - Use English in `Thought` part.
1174 - Write a small plan and finally summarize your next action (with its target element) in one
1175 sentence in `Thought` part.
1176
1177 ## User Instruction
1178 {TASK_INSTRUCTION}
1179
1180 {CURRENT_OBSERVATION}

```

High-Level Planner MLLM Prompt (Computer-Use: Ubuntu)

```

1176 You are an agent which follow my instruction and perform desktop computer tasks as instructed.
1177 You have good knowledge of computer and good internet connection and assume your code will run on
1178 a computer for controlling the mouse and keyboard.
1179 You are on Ubuntu operating system and the resolution of the screen is 1920x1080.
1180 For each step, you will get:
1181 - An observation of an image, which is the screenshot of the computer screen and you will predict
1182 the action of the computer based on the image.
1183 - Access to the following class and methods to interact with the UI:
1184 class Agent:
1185
1186     def click(self, instruction: str, num_clicks: int = 1, button_type: str = 'left', hold_keys:
1187         List = []):
1188         '''Click on the element
1189         Args:
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3019
3020
3021
3022
3023
3024
3025
3026
3027
3028
3029
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3069
3070
3071
3072
3073
3074
3075
3076
3077
3078
3079
3079
3080
3081
3082
3083
3084
3085
3086
3087
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3179
3180
3181
3182
3183
3184
3185
3186
3187
3188
3189
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3219
3220
3221
3222
3223
3224
3225
3226
3227
3
```

```

1188     instruction:str, describe the element you want to interact with in detail including the
1189     visual description and function description. And make it clear and concise. For
1190     example you can describe what the element looks like, and what will be the expected
1191     result when you interact with it.
1192     num_clicks:int, number of times to click the element
1193     button_type:str, which mouse button to press can be "left", "middle", or "right"
1194     hold_keys:List, list of keys to hold while clicking
1195     ...
1196
1197     def done(self, return_value: Union[Dict, str, List, Tuple, int, float, bool, NoneType] = None):
1198         '''End the current task with a success and the required return value'''
1199
1200     def drag_and_drop(self, starting_description: str, ending_description: str, hold_keys: List = []
1201         []:
1202         '''Drag from the starting description to the ending description
1203         Args:
1204             starting_description:str, a very detailed description of where to start the drag
1205             action. This description should be at least a full sentence. And make it clear and
1206             concise.
1207             ending_description:str, a very detailed description of where to end the drag action.
1208             This description should be at least a full sentence. And make it clear and concise.
1209             hold_keys:List list of keys to hold while dragging
1210             ...
1211
1212     def fail(self):
1213         '''End the current task with a failure, and replan the whole task.'''
1214
1215     def highlight_text_span(self, starting_phrase: str, ending_phrase: str):
1216         '''Highlight a text span between a provided starting phrase and ending phrase. Use this to
1217         highlight words, lines, and paragraphs.
1218         Args:
1219             starting_phrase:str, the phrase that denotes the start of the text span you want to
1220             highlight. If you only want to highlight one word, just pass in that single word.
1221             ending_phrase:str, the phrase that denotes the end of the text span you want to
1222             highlight. If you only want to highlight one word, just pass in that single word.
1223             ...
1224
1225     def hold_and_press(self, hold_keys: List, press_keys: List):
1226         '''Hold a list of keys and press a list of keys
1227         Args:
1228             hold_keys:List, list of keys to hold
1229             press_keys:List, list of keys to press in a sequence
1230             ...
1231
1232     def hotkey(self, keys: List):
1233         '''Press a hotkey combination
1234         Args:
1235             keys:List the keys to press in combination in a list format (e.g. ['ctrl', 'c'])
1236             ...
1237
1238     def open(self, app_or_filename: str):
1239         '''Open any application or file with name app_or_filename. Use this action to open
1240         applications or files on the desktop, do not open manually.
1241         Args:
1242             app_or_filename:str, the name of the application or filename to open
1243             ...
1244
1245     def scroll(self, instruction: str, clicks: int, shift: bool = False):
1246         '''Scroll the element in the specified direction
1247         Args:
1248             instruction:str, a very detailed description of which element to enter scroll in. This
1249             description should be at least a full sentence. And make it clear and concise.
1250             clicks:int, the number of clicks to scroll can be positive (up) or negative (down).
1251             shift:bool, whether to use shift+scroll for horizontal scrolling
1252             ...
1253
1254     def set_cell_values(self, cell_values: Dict[str, Any], app_name: str, sheet_name: str):
1255         '''Use this to set individual cell values in a spreadsheet. For example, setting A2 to "hello"
1256         would be done by passing {"A2": "hello"} as cell_values. The sheet must be opened before
1257         this command can be used.
1258         Args:
1259             cell_values: Dict[str, Any], A dictionary of cell values to set in the spreadsheet.
1260             The keys are the cell coordinates in the format "A1", "B2", etc.
1261             Supported value types include: float, int, string, bool, formulas.
1262             app_name: str, The name of the spreadsheet application. For example, "Some_sheet.xlsx".
1263             sheet_name: str, The name of the sheet in the spreadsheet. For example, "Sheet1".
1264             ...
1265
1266     def switch_applications(self, app_code):
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3020
3021
3022
3023
3024
3025
3026
3027
3028
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3070
3071
3072
3073
3074
3075
3076
3077
3078
3079
3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185
3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239
3240
3241
3242
3243
3244
3245
3246
3247
3248
3249
3250
3251
3252
3253
3254
3255
3256
3257
3258
3259
3260
3261
3262
3263
3264
3265
3266
3267
3268
3269
3270
3271
3272
3273
3274
3275
3276
3277
3278
3279
3280
3281
3282
3283
3284
3
```

```

1242     '''Switch to a different application that is already open
1243     Args:
1244         app_code:str the code name of the application to switch to from the provided list of
1245         open applications
1246         ...
1247
1248     def type(self, element_description: Optional[str] = None, text: str = '', overwrite: bool =
1249         False, enter: bool = False):
1250         '''Type text into a specific element
1251         Args:
1252             element_description:str, a detailed description of which element to enter text in.
1253             This description should be at least a full sentence.
1254             text:str, the text to type
1255             overwrite:bool, Assign it to True if the text should overwrite the existing text,
1256             otherwise assign it to False. Using this argument clears all text in an element.
1257             enter:bool, Assign it to True if the enter key should be pressed after typing the
1258             text, otherwise assign it to False.
1259             ...
1260
1261     def wait(self, time: float):
1262         '''Wait for a specified amount of time
1263         Args:
1264             time:float the amount of time to wait in seconds
1265             ...
1266
1267     The following rules are IMPORTANT:
1268     - If previous actions didn't achieve the expected result, do not repeat them, especially the last
1269     one. Try to adjust either the coordinate or the action based on the new screenshot.
1270     - Do not predict multiple clicks at once. Base each action on the current screenshot; do not
1271     predict actions for elements or events not yet visible in the screenshot.
1272     - You cannot complete the task by outputting text content in your response. You must use mouse and
1273     keyboard to interact with the computer. Call ```agent.fail()``` function when you think the task
1274     can not be done.
1275     - You must use only the available methods provided above to interact with the UI, do not invent
1276     new methods.
1277
1278     You should provide a detailed observation of the current computer state based on the full
1279     screenshot in detail in the "Observation:" section.
1280     Provide any information that is possibly relevant to achieving the task goal and any elements that
1281     may affect the task execution, such as pop-ups, notifications, error messages, loading states,
1282     etc..
1283     You MUST return the observation before the thought.
1284
1285     You should think step by step and provide a detailed thought process before generating the next
1286     action:
1287     Thought:
1288     - Step by Step Progress Assessment:
1289         - Analyze completed task parts and their contribution to the overall goal
1290         - Reflect on potential errors, unexpected results, or obstacles
1291         - If previous action was incorrect, predict a logical recovery step
1292     - Next Action Analysis:
1293         - List possible next actions based on current state
1294         - Evaluate options considering current state and previous actions
1295         - Propose most logical next action
1296         - Anticipate consequences of the proposed action
1297     Your thought should be returned in "Thought:" section. You MUST return the thought before the code.
1298
1299     You are required to use `agent` class methods to perform the action grounded to the observation.
1300     Return exactly ONE line of python code to perform the action each time. At each step (example:
1301     ```agent.click('Click \"Yes, I trust the authors\" button', 1, 'left')\n```)
1302     Remember you should only return ONE line of code, DO NOT RETURN more. You should return the code
1303     inside a code block, like this:
1304     ```python
1305     agent.click('Click \"Yes, I trust the authors\" button', 1, "left")
1306     ```
1307
1308     For your reference, you have maximum of {MAX_STEPS} steps, and current step is {CURRENT_STEP} out
1309     of {MAX_STEPS}.
1310     If you are in the last step, you should return ```agent.done()``` or ```agent.fail()``` according
1311     to the result.
1312
1313     Here are some guidelines for you:
1314     1. Remember to generate the corresponding instruction to the code before a # in a comment and only
1315     return ONE line of code.
1316     2. `agent.click` can have multiple clicks. For example, `agent.click('Click \"Yes, I trust the
1317     authors\" button', 2, "left")` is double click.
1318     3. Return ```agent.done()``` in the code block when you think the task is done (Be careful when
1319     evaluating whether the task has been successfully completed). Return ```agent.fail()``` in the
1320     code block when you think the task can not be done.
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678
2679
2680
2681
2682
2683
2684
2685
2686
2687
2688
2689
2690
2691
2692
2693
2694
2695
2696
2697
2698
2699
2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851
2852
2853
2854
2855
2856
2857
2858
2859
2860
2861
2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969
2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
2992
2993
2994
2995
2996
2997
2998
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3020
3021
3022
3023
3024
3025
3026
3027
3028
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3070
3071
3072
3073
3074
3075
3076
3077
3078
3079
3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185
3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212
3213
3214
3215
3216
3217
3218
3219
3220
3221
3222
3223
3224
3225
3226
3227
3228
3229
3230
3231
3232
3233
3234
3235
3236
3237
3238
3239
3240
3241
3242
3243
3244
3245
3246
3247
3248
3249
3250
3251
3252
3253
3254
3255
3256
3257
3258
3259
3260
3261
3262
3263
3264
3265
3266
3267
3268
3269
3270
3271
3272
3273
3274
3275
3276
3277
3278
3279
3280
3281
3282
3283
3284
3285
3286
3287
3288
3289
3290
3291
3292
3293
3294
3295
3296
3297
3298
3299
3300
3301
3302
3303
3304
3305
3306
3307
3308
3309
3310
3311
3312
3313
3314
3315
3316
3317
3318
3319
3320
3321
3322
3323
3324
3325
3326
3327
3328
3329
3330
3331
3332
3333
3334
```

1296 4. Whenever possible, your grounded action should use hot-keys with the `agent.hotkey()` action
 1297 instead of clicking or dragging.
 1298 5. Save modified files before returning `agent.done()`. When you finish modifying a file,
 1299 always save it before proceeding using `agent.hotkey(['ctrl', 's'])` or equivalent. Tasks may
 1300 involve multiple files. Save each after finishing modification.
 1301 6. If you meet "Authentication required" prompt, you can continue to click "Cancel" to close it.

1302 My computer's password is 'password', feel free to use it when you need sudo rights.

1303 Task Instruction:
 1304 `{TASK_INSTRUCTION}`

1305 Current Observation:
 1306 `{OBSERVATION}`

1307

1308 **Low-Level Planner MLLM Prompt - GTA1-7B (Computer-Use: Ubuntu)**

1309 You are an expert UI element locator. Given a GUI image and a user's element description, provide
 1310 the coordinates of the specified element as a single (x, y) point. The image resolution is height
 1311 `{HEIGHT}` and width `{WIDTH}`. For elements with area, return the center point.

1312 Output the coordinate pair exactly:
 1313 (x, y)

1314 Task Instruction: `{TASK_INSTRUCTION}`
 1315 `{CURRENT_OBSERVATION}`

1316

1317 **D.4 TASK VERIFIER**

1318 To verify if a task executed by the data collection policy is executed successfully, we employ a MLLM
 1319 as a task verifier that evaluates trajectories. The verifier is an MLLM that takes as input the task
 1320 instruction and the executed trajectory (represented as interleaved images and actions). For each
 1321 example it outputs the following in order: (i) screen_details: summarizing what the agent is doing
 1322 in the trajectory, (ii) reasoning: producing a Chain-of-Thought (Wei et al., 2023) justification of
 1323 whether the instruction has been completed, and (iii) result: issuing a final judgment of "success" or
 1324 "failure." We use GPT-4o (OpenAI et al., 2024) as the task verifier on expert-collected trajectories
 1325 and Qwen2.5-VL 32B Instruct](Bai et al., 2025) for RL training. The prompt used by both verifier
 1326 models is specified in Appendix D.4.

1327

1328 **Task Verifier Prompt (Mobile-Use: Android and Computer-Use: Ubuntu)**

1329 You are an AI assistant designed to help users evaluate whether a specific interaction with mobile
 1330 application was successful or not. At the beginning of the task you will be provided:

1331 1. a description of the task
 1332 2. a sequence of UI screenshots interleaved with actions taken in order to complete the specified
 1333 task

1334 For each action taken, you will be provided the action type, location shown on the screen with red
 1335 dot indicating where the interaction executed, and reason for the action taken.

1336 Your task is to carefully look at all the screenshots shown, the description of the task, and
 1337 actions taken to evaluate whether the specified task was completed or not. A task is only
 1338 considered completed if the screenshots only show the intended task being achieved. If the screens
 1339 show any unintended changes to the state of the app then the demonstration of task is treated as a
 1340 failure.

1341

Task instruction: `{TASK_INSTRUCTION}`

1342

Observations: `{OBSERVATIONS_WITH_ACTIONS}`

1343

You should output the evaluation response in the following JSON format:

{

1344

1345 "screen_details": "A bulleted list of all the changes on the screens as a result of all actions
 1346 executed. Even if the screen change do not correspond to the intended actions, you should still
 1347 describe the screen changes. Create a detailed list of all the UI changes that occurred on the
 1348 screen as a result of all the actions."

1349

"reasoning": "Step-by-step reasoning for the assesment of whether the agent was successful for
 the current step or if the agent has failed for the current step.",

1350	1351	Task Category	Description
1352	1353	1) create	creates a new entity in the app.
1354	1355	2) edit	modifies an existing entity in the app.
1356	1357	3) delete	deletes an entity from the app.
1358	1359	4) search	searches for something using the search bar. You can optionally filter the results based on the given criteria.
1360	1361	5) filter	filters the results based on the given criteria. You can optionally compose multiple filters to form a single filter. A filter can be direct, or it needs to be inferred using multi hop reasoning.
1362	1363	6) repeat operation	repeatedly applies a certain create/edit/delete/search operation multiple times to execute a task. For example, deleting multiple events in calendar or all events on a single day.
1364	1365	7) information retrieval	retrieves information from the app. This can involve answering questions about certain setting or data from an app
1366	1367	8) multi hop reasoning	involves reasoning across multiple steps or pieces of information to arrive at a conclusion. For example, answering a question that requires searching for data using multi-hop reasoning.
1368	1369	9) edit system setting	modifies a setting at the system level on the device.
1370	1371	10) edit app setting	modifies a setting within a specific application.
1372	1373	11) fill form	completes a form with the necessary information.
1374	1375	12) text editing	makes changes to text, such as editing, formatting, or restructuring.
1376	1377	13) cross app interaction	involves interacting with multiple applications to complete a task.
1378	1379	14) complex ui interaction	involves interacting with complex UI elements like timers and date time picker wheels.
1380	1381	15) composition	involves executing a task that requires completing multiple sub-tasks of different task categories to complete a longer horizon task. For example, tasks that require adding a new contact and deleting one involves executing both create and delete operation.

Table 7. Task Categories used for categorization of mobile-use tasks on Android and AndroidWorld.

E TASK CATEGORIZATION FOR ANALYSIS

To support the analysis in Sec. 4.2, each task is classified to multiple task categories that describe the types of skills required by an agent to complete the task. To annotate each task, we use an LLM (GPT-4o (OpenAI et al., 2024) in this case) and prompt it using the prompt specified in Appendix E using platform-specific categories specified in Tab. 7 for mobile-use agent tasks and in Tab. 8 for computer-use agent tasks.

Task Classifier Prompt (Mobile-Use: Android and Computer-Use: Ubuntu)

You are a helpful assistant that can help me analyze the task can identify the task categories which uniquely describe the types of tasks and skills required by a GUI agent to complete the task. As input you will be provided a list of task instructions and your task is to output the list of task categories that apply to the specified instructions. Next, we will list the type of task categories you need to label each task using following categories:

{TASK_CATEGORIES}

For each task instruction, you have to output a response in JSON format.

- A list of task categories that apply to the specific task instruction.
- Use only the list of task categories listed above as part of the response.

1404	Task Category	Description
1405	1) create	creates a new entity in the application.
1406	2) edit	modifies an existing entity in the application.
1407	3) delete	deletes an entity from the application.
1408	4) filter	filters search results or data to operate on based on the given criteria. You can optionally compose multiple filters to form a single filter. A filter can be direct, or it needs to be inferred using multi hop reasoning.
1409	5) repeat operation	repeatedly applies a certain create/edit/delete/search operation multiple times to execute a task. For example, deleting multiple events in calendar or all events on a single day.
1410	6) information retrieval	retrieves information from the app. This can involve answering questions about certain setting or data from an app
1411	7) multi hop reasoning	involves reasoning across multiple steps or pieces of information to arrive at a conclusion. For example, answering a question that requires searching for data using multi hop reasoning.
1412	8) file management	involves managing files and folders on the computer.
1413	9) coding	involves writing code to accomplish a specific task or solve a problem.
1414	10) web search	searches for information on the web using a search engine.
1415	11) device search	searches for information on the computer file system.
1416	12) cross app tasks	involves interacting with desktop multiple applications to complete a task.
1417	13) edit app setting	modifies a setting within a specific application.
1418	14) edit system settings	modifies a setting at the system level on the computer.
1419	15) composition	involves executing a task that requires completing multiple sub tasks of different task categories to complete a longer horizon task. For example, tasks that require adding a new contact and deleting one involves executing both create and delete operation.
1420	16) content editing	tasks that require editing content such as text, images, videos, or audio.
1421	17) spreadsheet editing	tasks that require editing or manipulating data in a spreadsheet application.
1422	18) document editing	tasks that require editing or formatting documents, such as word processing files or PDFs.
1423	19) bash scripting	tasks that require writing or executing bash scripts or commands using the terminal or code editors.

Table 8. Task Categories used for categorization of computer-use tasks on Ubuntu and OSWorld.

1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

```

1458
1459     Here are some examples of task categorization:
1460     Instructions: [
1461         "Create a timer with {hours} hours, {minutes} minutes, and {seconds} seconds. Do not start the
1462         timer.",
1463         "Add the expenses from expenses.jpg in Simple Gallery Pro to pro expense.",
1464         "Change the theme to {theme} in the Audio Recorder app."
1465     ]
1466     Output:
1467
1468     [
1469         {
1470             "task_instruction": "Create a timer with {hours} hours, {minutes} minutes, and {seconds} seconds. Do not start the timer.",
1471             "task_categories": [
1472                 "create",
1473                 "complex ui interaction"
1474             ]
1475         },
1476         {
1477             "task_instruction": "Add the expenses from expenses.jpg in Simple Gallery Pro to pro expense.",
1478             "task_categories": [
1479                 "create",
1480                 "cross app interaction",
1481                 "fill form"
1482             ]
1483         },
1484         {
1485             "task_instruction": "Change the theme to {theme} in the Audio Recorder app.",
1486             "task_categories": [
1487                 "edit app setting"
1488             ]
1489         }
1490     ]
1491
1492     Next, you will be provided 10 tasks as input and your task is to output a JSON in specified format:
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

```

F TASK PROPOSER ANALYSIS

In this section, we describe the prompts and setup used to synthesize data used for *No Exploration* and *Iterative Exploration* task generators for comparison with AUTOPLAY.

No Exploration: We use the prompt mentioned in Appendix F to instantiate the *No Exploration* method. In order to provide high-quality environment context we manually write text descriptions for all 20 android applications we use for generating training data. Each of these description aims to enumerate the list of functionalities accessible in an application. This information is gathered by us manually by exploring each android application manually which is a key limitation of this method. We show example descriptions of two applications in Appendix F.

Iterative Exploration: To instantiate the *Iterative Exploration* method we use the public implementation of AgentSynth (Xie et al., 2025) as it demonstrates a effective pipeline for iterative exploration for data collection for web agents. We adapt AgentSynth (Xie et al., 2025) for mobile-use domain by using the data collection policy described in Sec. 3.3 as the task executor. We run the iterative task proposal and execution for k turns where k is randomly sampled to be between 3 – 8 turns and each turn can run up to 7 steps. We use 7 steps as a limit for each subtask as we qualitatively find that subtasks proposed by such methods can be completed with 4 – 7 steps on average. After all k turns are executed the complete trajectory is hindsight relabelled using descriptions of each subtask and whether they were successfully completed or not. Using this setup, we run this method to generate $5k$ tasks and use trajectories where atleast 50% subtasks succeeded as our training dataset.

No Exploration Task Generator MLLM Prompt

```

1505
1506
1507     You are a capable UI understanding agent. You will be provided a app description for the Android
1508     app. Your task is to propose a list of tasks that a useful UI understanding agent should be able
1509     to complete. Use the app description provided to ground these tasks to the ones that are feasible.
1510     Propose as many diverse tasks as possible to cover broad range of features.
1511

```

```

1505
1506
1507     App description: {app_description}
1508
1509
1510
1511

```

1513 For all described functionalities output a list of up to {NUM_TASKS} unique tasks that can be
1514 executed in the app as a JSON which should be in the following format:

```
1515 tasks = [
1516     {
1517         "thought": "<Detailed thought and reasoning for the proposed task, why is the task simple
1518         enough to execute and whether it satisfies the task proposal guidelines>",
1519         "instruction": "<natural language instruction describing a task with/without template
1520         params with name of the app>,
1521         "tag": <few words describing the type of task>,
1522         "app_name": <name of the app>,
1523         "template_params": {
1524             "param_name": {
1525                 "description": <param description>,
1526                 "possible_values": [<list of 5 random values>]
1527             }
1528         }
1529     }
1530 ]
1531
1532 For each functionality/feature make sure to follow these guidelines while generating tasks:
1533 1. Create all possible tasks. For example, if there is a clock app on the phone which has
1534 worldwide clock feature then you can create tasks like: "Add worldwide clock for <city>", "Remove
1535 worldwide clock for <city>", "Add worldwide clock for <city_1, city_2> and reorder to put <city_2>
1536 before <city_1>". Make sure to add the <param> details in "template params" in JSON to allow
1537 creation of diverse tasks.
1538 2. In any instruction if there are going to be templated parameters then add it in following
1539 format: {param_name}
1540 3. If a task requires creating/adding/editing/deleting any information/event/entry make sure all
1541 required entities for completing the task are parameterized/templated. For example:
1542     a.) For a task that requires creating notes both the name of the new note and content of the
1543 note should be specified as parameters.
1544     b.) Similarly, if a task requires editing some details describe the exact edit you'd like to
1545 make, what entity needs to be edited. When proposing such tasks, only ask edits to existing
1546 entities shown in the images from the app for such tasks. If no such entities exist ask the
1547 agent to first create one and then edit it.
1548     c.) If a task requires deleting an entity, describe the entity that needs to be deleted. When
1549 proposing such a task, only ask to delete an entity that already exists as shown in the images
1550 from the app for such tasks. If no such entities exist ask the agent to first create one and then
1551 delete it.
1552     d.) If a task requires copying an entity, describe the entity that needs to be copied and
1553 where it should be copied to (ex: which folder or date). When proposing such a task, only ask
1554 to copy an entity that already exists as shown in the images from the app for such tasks. If
1555 no such entities exist ask the agent to first create one and then copy it.
1556     e.) For all parameters that are templated if any parameter refers to an entity for edit,
1557 delete or copy task make sure that the list of possible values only contain entities that
1558 exists or will be created.
1559 4. Ensure that instructions are unambiguous and clearly describes a actual task that can be
1560 performed on the app.
1561 5. Do not include tasks that require accessing/uploading/capturing content from real world. For
1562 example, scanning a document, recording a new video, taking a picture using camera.
1563 6. Task instructions should be natural user requests someone might actually ask a capable UI agent
1564 to complete on the app.
1565 7. If any supports browsing information from the web or library make sure to include that in the
1566 task instruction. For example, if a task requires searching for a book on the web, make sure to
1567 include that in the task instruction.
```

App Description Examples (Mobile Use: Android)

```
1554
1555     ### Audio Recorder app
1556
1557     The Audio Recorder app is a versatile and user-friendly application designed for recording audio
1558     with a range of customizable settings. It caters to personal, educational, and professional audio
1559     recording needs by offering options to select recording formats, sample rates, bitrates, and
1560     channel counts. The app also provides features for theme customization, file management, and
1561     sharing, making it suitable for various audio capture and management tasks.
1562
1563     "Theme Customization: Users can personalize the app's appearance by selecting from themes such as
1564     Blue Gray, Black, Teal, Blue, Purple, Pink, Orange, Red, and Brown. This feature enhances the user
1565     experience by allowing visual customization."
1566     "Recording Format Selection: Users can choose between M4a, Wav, and 3gp formats. M4a is
1567     recommended for its good quality and small size, Wav is uncompressed and takes more space, and 3gp
1568     is suitable for saving space. This feature allows users to tailor the recording quality and file
1569     size to their needs."
1570     "Sample Rate Selection: Offers options for 8kHz, 16kHz, 22kHz, 32kHz, 44.1kHz, and 48kHz, enabling
1571     users to select the desired audio quality. Higher sample rates provide better audio quality.",
```

1566 "Bitrate Selection: Users can choose from 48 kbps, 96 kbps, 128 kbps, 192 kbps, and 256 kbps to
 1567 control the audio quality and file size. Higher bitrates result in better audio quality.",
 1568 "Channel Count: Options for Stereo and Mono recording, providing flexibility in audio capture.
 1569 Stereo offers two channels for richer sound, while Mono uses a single channel.",
 1570 "File Management: Includes a file browser to access recorded files, options to rename, view
 1571 detailed information (format, bitrate, channel count, sample rate, duration, size, file location,
 1572 creation date), organize files by date, and a trash feature where deleted files are stored for 60
 1573 days before permanent deletion. This feature aids in efficient file organization and retrieval.",
 1574 "Recording Interface: Displays recording time, file size, format, and sample rate during
 1575 recording, with easy access to settings and file management. Users can start, pause, and stop
 1576 recordings, and view a visual waveform of the audio being recorded.",
 1577 ======
 1578 ======
 1579 **### Simple Calendar Pro**
 1580
 1581 This app is a versatile calendar application designed to help users efficiently manage their
 1582 schedules for both personal and professional use. It offers a range of features including event
 1583 and task management, multiple calendar views, import/export functionality, and customization
 1584 options. Users can create and organize events and tasks, set reminders, and customize their
 1585 calendar experience with color and notification settings. The app supports adding holidays,
 1586 birthdays, and anniversaries, and provides options for printing schedules and navigating through
 1587 dates.
 1588 "Calendar View: Displays various views such as daily, weekly, monthly, and yearly, allowing users
 1589 to navigate through months and view scheduled activities. Users can customize their view
 1590 preferences and highlight weekends.",
 1591 "Event Creation: Users can create new events with details such as title, location, description,
 1592 start and end time, reminders, repetition settings, event type, and color customization. Events
 1593 can be added, edited, and deleted.",
 1594 "Task Creation: Users can add tasks to their calendar, similar to events, to manage to-do lists
 1595 and deadlines. Tasks can be organized and viewed in different calendar views.",
 1596 "Search Functionality: Users can search for specific events or tasks using keywords, making it
 1597 easy to find events quickly.",
 1598 "Import/Export Events: Users can import events from an .ics file and export their calendar events
 1599 to an .ics file, facilitating easy sharing and backup of schedules.",
 1600 "Print Functionality: Users can print their schedules, with options to select the number of copies
 1601 and paper size. The app supports saving the schedule as a PDF or printing via connected printers.",
 1602 "Add Holidays and Contacts: Users can add holidays, contact birthdays, and anniversaries to their
 1603 calendar, ensuring they never miss important dates. The app may request access to contacts for
 1604 this feature.",
 1605 "Event Reminders: Users can set reminders for events to receive notifications. Options include
 1606 setting no reminder or adding new birthdays and anniversaries automatically.",
 1607 "Settings Customization: Offers options to customize colors, widget colors, language settings,
 1608 time format (24-hour), week start day, and weekend highlighting. Users can also customize
 1609 notifications, choose audio streams for reminders, and enable vibration for notifications.",
 1610 "Color Customization: Users can change the theme and app icon color, with a warning about
 1611 potential issues with some launchers.",
 1612 "Go to Date: Allows users to quickly navigate to a specific date in the calendar.",
 1613 "About: Provides information about the app."
 1614
 1615
 1616
 1617
 1618
 1619

G ACTION SPACE

1606 For mobile-use agents, we use the low-level action space from AndroidWorld (Rawles et al., 2024)
 1607 benchmark. The full list of actions used by our data collection policy and AUTOPLAY finetuned
 1608 models is listed in Tab. 9.

1609 For training desktop-use agents, we use a hybrid action space implemented in Agent-S Agashe et al.
 1610 (2025) and used by prior works (Yang et al., 2025). The full list of actions used by our data collection
 1611 policy and AUTOPLAY finetuned models is listed in Tab. 10. At the low-level these actions are
 1612 implemented using pyautogui APIs widely used in methods on the OSWorld benchmark (Xie et al.,
 1613 2024).

1620

1621

1622

1623

1624

1625

1626

1627

Action	Description
click	Click on (x, y) coordinate
long_press	Tap and hold on (x, y) coordinate
input_text	Types ‘text’ in a textbox if active
open_app	Opens android application given as text argument
scroll	Scrolls screen in direction (up or down)
wait	Waits for ‘time’ specified in seconds
navigate_back	Navigates back to previous screen
navigate_home	Navigates to home screen of the device
terminate	Ends the current episode

Table 9. Action Space used for mobile-use agent evaluated on AndroidWorld.

1638

1639

1640

1641

1642

1643

1644

1645

1646

1647

1648

1649

1650

1651

1652

1653

Action	Description
click	Click on (x, y) coordinate, num_click times, using either left or right click button, while holding keys given by ‘hold_keys’
hold_and_press	Holds list of keys ‘hold_keys’ and presses keys given by ‘press_keys’
hotkey	Invokes a hotkey combination given by ‘keys’
open	Opens the file name or application specified as the argument
scroll	Clicks on (x, y) coordinate, scrolls in up or down direction while holding shift key if ‘shift’ argument is set to true
set_cell_values	Sets individual cells in a spreadsheet to certain values
switch_applications	Switch to a different application that is already open specified as argument
type	Type specified text in an element
wait	Wait for a specified amount of time in seconds
fail	Ends the current episode and returns failure response
done	Ends the current episode

Table 10. Action Space used for computer-use agent evaluated on OSWorld. We use pyautogui action space to implement the specified actions.

1666

1669

1670

1671

1672

1673