

NATURALGAIA: PUSHING THE FRONTIERS OF GUI AGENTS WITH A CHALLENGING BENCHMARK AND HIGH-QUALITY TRAJECTORY DATASET

Anonymous authors

Paper under double-blind review

ABSTRACT

The rapid advancement of Large Language Model (LLM)-driven Graphical User Interface (GUI) agents is significantly hampered by the profound limitations of existing evaluation benchmarks in terms of accuracy, reproducibility, and scalability. To address this critical gap, we introduce NaturalGAIA, a novel benchmark engineered on the principle of Causal Pathways. This design paradigm structures complex tasks into a series of programmatically verifiable atomic steps, ensuring a rigorous, fully automated, and reproducible standard for assessment. Concurrently, to mitigate the inherent capability deficits of agents, we developed LightManus, a hierarchical agent architecture specifically optimized for long-horizon tasks. We leveraged this agent to generate a high-quality, human-verified trajectory dataset that uniquely captures diverse and even self-correcting interaction patterns of LLMs. We then utilized this dataset to perform Reinforcement Fine-Tuning (RFT) on the Qwen2.5-VL-7B model. Our experiments reveal that NaturalGAIA presents a formidable challenge to current state-of-the-art LLMs; even the top-performing Claude-sonnet-4 achieved a Weighted Pathway Success Rate (WPSR) of only 34.6%. Moreover, while RFT substantially improved the smaller model’s GUI execution capabilities (WPSR increased from 3.3% to 10.8%), its performance degraded sharply when handling complex scenarios. This outcome highlights the inherent capability ceiling of smaller models when faced with comprehensive tasks that integrate perception, decision-making, and execution. This research contributes a rigorous evaluation standard and a high-quality dataset to the community, aiming to guide the future development of GUI agents.

1 INTRODUCTION

The emergence of LLM-driven GUI agents (DeepSeek-AI et al., 2025; OpenAI et al., 2024b; Qwen et al., 2025) represents a significant advance over traditional automation (Yang et al., 2023; Zhang et al., 2023; Tang et al., 2025; Liu et al., 2025a; Durante et al., 2024), yet their progress is difficult to assess due to severe limitations in existing benchmarks. Current evaluation methods often rely on ambiguous final-state verification, which depends on fallible MLLMs (Bai et al., 2023; OpenAI et al., 2024a) and costly manual checks, thus compromising the scalability, reproducibility, and automation of evaluations (Mialon et al., 2023; Rawles et al., 2024). Moreover, these benchmarks inadequately measure performance on long-horizon tasks and fail to consider diverse valid solution paths, resulting in insufficient evaluation comprehensiveness (Liu et al., 2023; Li et al., 2025a; Chezelles et al., 2025).

Concurrently, these agents exhibit significant deficits in handling complex, real-world workflows, a problem exacerbated by prevailing learning paradigms (Koh et al., 2024b; Xu et al., 2024). Agents trained on successful execution traces Wang et al. (2025); Jiang et al. (2025); Qiu et al. (2025) tend to form rigid “shortcut” strategies. While effective for repeated tasks in static environments, this approach leads to poor generalization and robustness in novel scenarios, indicating a reliance on memory over transferable knowledge (Cheng et al., 2024; He et al., 2024; Gou et al., 2025). These generalization failures fundamentally stem from the LLM’s deficiencies in GUI operational knowledge and long-range planning (Ye et al., 2025b; Muennighoff et al., 2025; Huang et al., 2024), a constraint rooted in the scarcity of high-quality, structured trajectory data He et al. (2025).

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

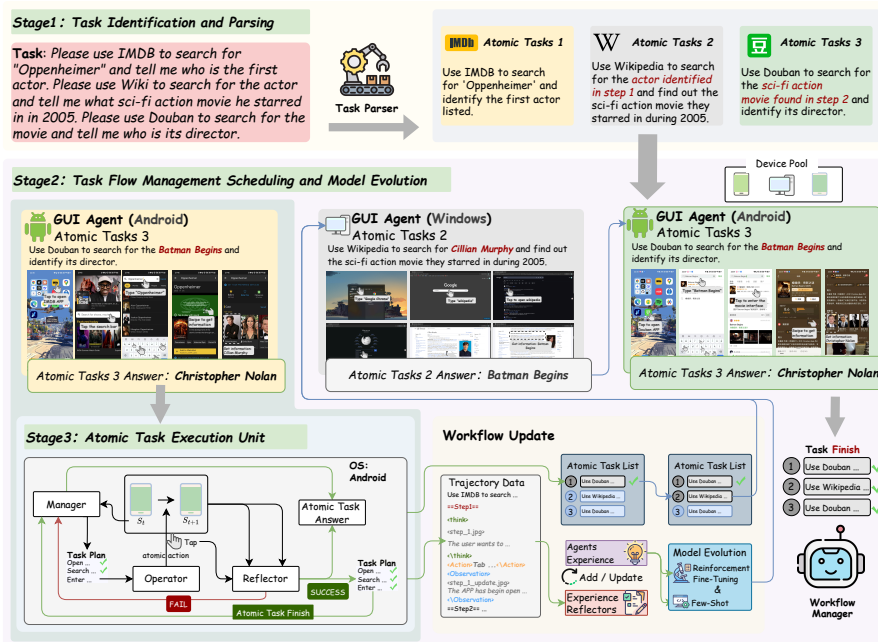


Figure 1: An overview of our three-stage framework **LightManus**, executing a complex, cross-platform task from NaturalGAIA. **Task Parsing** (Stage 1): A complex instruction is decomposed into a sequence of atomic tasks that form a Causal Pathway. **Workflow Management & Evolution** (Stage 2): The system schedules these tasks across different devices and uses collected Trajectory Data for model evolution via RFT. **Atomic Task Execution Unit** (Stage 3): This core unit, comprising a Manager, Operator, and Reflector, handles the execution of each atomic task.

To address these issues, we introduce NaturalGAIA, a novel benchmark with 276 tasks designed to address critical evaluation challenges in GUI agents. Unlike benchmarks such as GAIAMialon et al. (2023) that focus on flexible, goal-oriented reasoning, NaturalGAIA is engineered to rigorously test an agent’s ability to execute complex, structured procedures reliably. This is achieved through Causal Pathways (CPs)—task structures that model realistic, high-frequency user workflows across multiple platforms, thereby simulating how humans sequentially use tools to achieve goals. Each pathway is composed of interdependent, programmatically verifiable atomic steps. An error at any step triggers a deterministic “Path Collapse”, which not only provides an unambiguous failure signal but also pinpoints the source of the failure. Consequently, NaturalGAIA ensures task complexity and generalizability while significantly enhancing evaluation timeliness, accuracy, and automation, thereby establishing a rigorous and reproducible standard for assessing GUI agents.

Furthermore, advanced learning paradigms were explored to advance the capabilities of GUI agents. We developed a hierarchical agent architecture featuring a Task Parser and a Workflow Manager to improve long-horizon task handling through effective decomposition and dynamic scheduling. Additionally, we constructed a high-quality GUI operation trajectory dataset. Generated by our agent and iteratively refined through expert validation, this dataset ultimately contains 523 trajectories with detailed planning rationale. It was subsequently used to investigate the efficacy of Reinforcement Fine-Tuning (RFT) and serves as a foundational resource for developing next-generation agents. A representative example is presented in Figure 1.

Experimental results demonstrate that our hierarchical agent, LightManus (driven by Claude-sonnet-4), significantly outperforms mainstream baselines with a 42.9% Pass@1 success rate, yet achieves only a 34.6% Weighted Pathway Success Rate (WPSR) on our NaturalGAIA. We also validated the efficacy of our collected data through RFT, which substantially enhanced the Qwen2.5-VL-7B model’s capability by increasing its WPSR from 3.3% to 10.8%. However, this enhancement does not generalize to complex, long-horizon tasks where performance degrades sharply. These findings

108 reveal a critical gap between specialized, data-driven skill acquisition and the generalizable, long-
109 range planning required for true GUI autonomy.

110 The main contributions of this paper can be summarized as follows:

- 111
- 112 • We introduce NaturalGAIA, a novel benchmark meticulously designed around Causal Path-
113 ways principles and equipped with a multi-level evaluation framework to precisely evaluate
114 agent performance on complex, long-horizon tasks. Experimental results reveal that even
115 state-of-the-art models encounter a formidable challenge on this suite.
- 116 • We construct and publicly release a high-quality dataset of GUI operation trajectories, com-
117 prising 523 trajectories curated through an iterative process of LLM generation and rigor-
118 ous human verification, to provide a critical resource for future agent training.
- 119 • We conduct an in-depth investigation into trajectory-based reinforcement fine-tuning. This
120 study offers initial evidence and key insights into leveraging this method efficiently to en-
121 hance agent capabilities.
- 122

123 2 RELATED WORK

124 The advancement of GUI agents faces three principal challenges. First, an evaluation-realism
125 dilemma persists in benchmarks. While environments have evolved towards higher fidelity, from
126 web-based Zhou et al. (2023); Koh et al. (2024a) to mobile-centric Chai et al. (2025) and long-
127 horizon tasks Ye et al. (2025a), this pursuit of realism compromises scientific rigor. The dynamic
128 nature of modern applications impedes reproducibility Garg et al. (2025), and the existence of mul-
129 tiple valid solution paths renders fixed trajectory matching unreliable Zhou et al. (2023); Liu et al.
130 (2023), often necessitating costly or fallible evaluations Chai et al. (2025); Mialon et al. (2023). Sec-
131 ond, a data bottleneck constrains the development of increasingly sophisticated agents. Advances
132 in high-resolution visual perception Hong et al. (2024); Li et al. (2025b); Cheng et al. (2024) and
133 modular architectures Zhang et al. (2025; 2024); Wang et al. (2024) demand structured, fine-grained
134 data that existing datasets lack. Third, a deficit of experiential data hinders the shift from simple im-
135 itation learning Rasheed et al. (2024); Deng et al. (2023) to more advanced paradigms like learning
136 from self-exploration Wang et al. (2023); Jiang et al. (2025) or trial-and-error Shinn et al. (2023);
137 Qiu et al. (2025); Wang et al. (2024), which require diverse trajectory data. This work addresses
138 these gaps directly.

139 NaturalGAIA resolves the evaluation dilemma by providing a stable, reproducible, and fully auto-
140 mated benchmark with programmatically verifiable answers. Concurrently, our trajectory dataset
141 addresses both the data bottleneck and the experiential data deficit by supplying structured and di-
142 verse interaction data to foster the development of more capable and robust agents.

143 3 NATURALGAIA

144 While prominent benchmarks such as GAIA excel at assessing flexible, goal-oriented reasoning,
145 they lack rigorous testing of an agent’s ability to reliably execute complex, structured procedures
146 grounded in natural user tasks. To fill this gap, NaturalGAIA introduces Causal Pathways (CPs)—a
147 task architecture of strict, causally-dependent sequences explicitly designed to model these authen-
148 tic, cross-platform procedures. In this framework, any procedural error triggers a deterministic Path
149 Collapse, yielding an unambiguous, process-level failure signal that transcends the limitations of
150 final-state evaluations. NaturalGAIA thus functions as a precise diagnostic tool, designed to probe
151 an agent’s core capabilities in long-range planning, memory, and operational fidelity.

152 3.1 ATOMIC TASKS: THE FUNDAMENTAL NODES OF CAUSAL PATHWAYS

153 At the core of NaturalGAIA are *atomic tasks*—the smallest indivisible units of interaction. Each
154 atomic task is designed to retrieve a single, specific piece of information (i.e., the “answer”) through
155 a sequence of primitive GUI operations. Importantly, every atomic task is associated with a unique
156 and deterministic correct answer, allowing for objective, automated evaluation. This granularity en-
157 ables the benchmark to assess an agent’s performance in perception, interface manipulation, and rea-
158 soning independently. By isolating each task’s input, action sequence, and output, NaturalGAIA fa-
159
160
161

162 cilitates precise error analysis and supports targeted improvements in sub-capabilities. These atomic
 163 tasks serve as the fundamental nodes that constitute each Causal Pathway.

165 3.2 CAUSAL PATHWAYS: STRUCTURING TASKS WITH LOGICAL DEPENDENCIES

166
 167 The core construct of NaturalGAIA is the Causal Pathway (CP), an executable representation of pro-
 168 cedural knowledge for structured tasks. A key design principle is the decoupling of a CP’s abstract
 169 logic from its physical execution across diverse platforms, such as different applications or web-
 170 sites. This allows for the flexible generation of varied task instances to mitigate data contamination
 171 and model overfitting, while simultaneously upholding the rigorous evaluation of underlying causal
 172 dependencies.

173 Nodes within a CP are linked by strict causal dependencies, where an error at any step triggers
 174 a “Path Collapse”—a deterministic failure of the entire sequence. This mechanism provides an
 175 unambiguous signal, pinpointing the exact moment an agent’s execution deviates from the required
 176 causal logic. Consequently, NaturalGAIA evaluates not just task completion but, more critically, an
 177 agent’s capacity for long-range planning, abstract reasoning, and robust execution.

178 The specific construction principles and processes of CPs and tasks in NaturalGAIA are shown in
 179 Appendix D.1 and D.2 respectively.

181 3.3 MULTI-LEVEL EVALUATION FRAMEWORK

182
 183 To precisely evaluate agent performance within NaturalGAIA and meticulously identify the root
 184 causes of errors, this study constructs a hierarchical evaluation framework. This framework en-
 185 ables comprehensive scrutiny from macroscopic, difficulty-weighted task completion down to mi-
 186 croscopic error attribution.

187
 188 **Level 1: Weighted Pathway Success Rate (WPSR)** To account for task complexity, we introduce
 189 a difficulty score $D_{j,i}$ for each task instance i of a Causal Pathway CP_j . This score is proportional
 190 to the number of nodes in the pathway and the number of distinct applications involved in its cor-
 191 responding task. WPSR is defined by weighting each successful completion ($\mathbb{S}_{\text{task}}(j, i) = 1$) by
 192 its normalized difficulty $w_{j,i} = D_{j,i} / \sum_{k,l} D_{k,l}$, where the sum in the denominator is over all task
 193 instances across all pathways.

$$194 \text{WPSR} = \sum_{j,i} w_{j,i} \cdot \mathbb{S}_{\text{task}}(j, i)$$

195
 196 WPSR serves as a holistic metric for an agent’s final task completion capability, weighted by diffi-
 197 culty.

198
 199 **Level 2: Fine-grained Traversal Metrics** To quantify partial progress and resilience, we intro-
 200 duce two complementary metrics that assess traversal quality.

201 **Mean Atomic Tasks Completion Ratio (MATCR):** This metric quantifies an agent’s ability to
 202 successfully complete sequences of atomic tasks. For each task sequence j within the benchmark,
 203 a completion ratio $R_j = k_j/n_j$ is calculated, where k_j is the number of consecutive atomic tasks
 204 successfully executed from the start, and n_j is the total number of atomic tasks in sequence j .

$$205 \text{MATCR} = \frac{1}{N} \sum_{j=1}^N \frac{k_j}{n_j}$$

206
 207
 208
 209 **Positional-Weighted Atomic Tasks Success Rate (p-ATSR):** This metric evaluates an agent’s abil-
 210 ity to maintain long-term coherence by assigning greater weight to successes in the later stages of
 211 a Causal Pathway. Let n_j be the total number of atomic tasks in pathway CP_j . We introduce a
 212 positional weight, denoted as $p(i)$, which is a monotonically increasing function of the step index i .
 213 The p-ATSR is then defined as:

$$214 \text{p-ATSR} = \frac{\sum_{j=1}^N \sum_{i=1}^{n_j} p(i) \cdot \mathbb{S}_{\text{atomic}}(j, i)}{\sum_{j=1}^N \sum_{i=1}^{n_j} p(i)}$$

where $\mathbb{S}_{\text{atomic}}(j, i) = 1$ for a success at step i of pathway j , and 0 otherwise.

In summary, MATCR assesses an agent’s foundational reliability by quantifying its average execution length, whereas p-ATSR places greater emphasis on its long-term coherence by rewarding success in the later stages of a task.

Level 3: Error Attribution Analysis For atomic tasks that fail, as identified by the Level 2 metrics, this stage involves analyzing the atomic operation sequence to attribute Execution Errors (EE) to one of three primary types:

- *Knowledge Deficit (KD)*: The agent lacks the necessary domain-specific or procedural knowledge regarding application operations or function invocation required to achieve the objective.
- *Perceptual Error (PE)*: Although the agent possesses the correct intent and planning, it fails during the information extraction phase due to misreading or misinterpretation of visual or textual information from the interface.
- *Operational Error (OE)*: The agent understands the task and correctly perceives the interface, but fails to take appropriate or precise actions during specific operations.

3.4 TASK DIFFICULTY STRATIFICATION

To enable a systematic evaluation of agent capabilities, tasks in NaturalGAIA are stratified into three difficulty levels. This stratification is defined by the structural properties of their underlying Causal Pathways, including their length (number of nodes), the complexity of transitions between nodes, and the diversity of applications involved.

Level 1: Consisting of short Causal Pathways involving 1-2 atomic tasks (nodes). These tasks typically operate within 1-2 applications and assess an agent’s basic sequential execution and information-passing capabilities.

Level 2: Featuring Causal Pathways of intermediate length, containing 3-4 nodes. These pathways often require transitions across 3-5 different applications or involve complex interactions within a smaller set of applications, testing an agent’s ability to maintain context over a longer duration.

Level 3: Involving complex and lengthy Causal Pathways spanning 5-7 nodes and 3-7 applications. Successfully traversing these pathways demands advanced planning, long-term memory, and robust handling of intricate inter-node dependencies, posing a significant challenge to the agent’s cognitive and operational architecture.

Finally, we constructed 276 tasks. Data statistics and examples of tasks at different levels are shown in the Appendix D.3 and D.4.

4 METHOD

4.1 CROSS-PLATFORM AGENT ARCHITECTURE

To facilitate the generation of high-quality, long-horizon task trajectories, this study introduces LightManus, a specialized framework engineered for complex, cross-platform operations. Its hierarchical architecture is systematically organized into three distinct stages. The first stage, Task Parsing, analyzes the high-level user instruction and decomposes it into a logical sequence of executable sub-tasks. The second stage, the Workflow Manager, orchestrates this sequence by dynamically scheduling and delegating each sub-task to the appropriate device. The final stage is the Atomic Task Execution Unit, which executes the ground-level operations. This unit consists of specialized low-level agents—Mobile-Agent-e Wang et al. (2025) for Android and PC-Agent Liu et al. (2025b) for PC—that handle direct interaction with the graphical user interface. This structured, three-stage design provides a robust mechanism for generating high-fidelity execution trajectories and optimizes long-horizon task performance by ensuring seamless coordination across disparate operating systems.

4.2 HUMAN-VERIFIED TRAJECTORY COLLECTION

To construct a high-quality dataset aimed at enhancing the capabilities of GUI agents, we implemented a **rigorous** two-phase data collection pipeline. This methodology is designed to synergize the scalability and diverse exploratory power of a Large Language Model (LLM)-driven agent with the precision and reliability of human experts. The core motivation is to create a dataset that not only ensures procedural correctness but also captures a rich diversity of interaction behaviors, thereby transcending the traditional paradigm of singular, “optimal” paths.

Phase 1: LLM-Driven Trajectory Generation The data generation process was initiated using 265 unique atomic tasks as prompts. Our agent, LightManus, powered by Gemini-2.5-Pro, was tasked with autonomously completing these tasks. For each task, LightManus engaged in a **cyclical action process** at every step, which can be decomposed as follows:

- **Perception and State Understanding:** The agent first captures and analyzes the current user interface (UI) screenshot.
- **Deliberation and Planning:** The agent integrates the task objective with visual information to perform multi-level reasoning. This includes a high-level `Thought` process to maintain the overall plan and a micro-level `Action Thought` to determine the specific UI operation to be executed.
- **Action Execution:** The agent formulates and dispatches the determined action command.
- **Observation and Reflection:** Following the action, the agent captures and observes the new interface state, then provides feedback on the outcome to reflect on its progress toward the goal.

This automated procedure yielded an initial corpus of over 600 successfully executed trajectories. However, this raw corpus suffered from a critical flaw that necessitated curation: in certain instances, the model would bypass genuine UI interaction by directly generating answers from its internal, parametric knowledge. This highlighted the indispensable value of the human verification phase.

Phase 2: Expert Verification and Curation Every successfully generated trajectory was subjected to a meticulous review by a team of two graduate researchers. This review involved a detailed source-trace analysis of all execution information, including step-by-step screenshots and corresponding log files. Through an iterative refinement of this analysis, 523 trajectories were ultimately curated for the final dataset, adhering to the following principles:

- **Enforcing Grounded Execution:** A trajectory was deemed valid only if all of the agent’s actions were strictly grounded in the visual information presented on-screen. This principle explicitly filters out any solutions derived from the model’s internal knowledge, ensuring the data reflects genuine agent-environment interaction.
- **Preserving Behavioral Diversity and Corrective Processes:** Our curation criteria deliberately deviate from selecting only the “most efficient” paths, aiming instead to retain a variety of valid interaction strategies. Furthermore, we intentionally preserved trajectories containing recoverable *error steps*—cases where the agent successfully self-corrects through reflection after executing an erroneous action or exploring a wrong path. We posit that these trajectories, containing corrective processes, provide an invaluable learning signal for training more robust agents capable of reasoning about and recovering from their own errors.

4.3 FINAL DATA STRUCTURE

The curated dataset is composed of structured data points, each containing the initial instruction, the final solution, and the interaction trajectory. A complete trajectory T is represented as a sequence of time steps, $T = (t_1, t_2, \dots, t_n)$. Each time step t_i is defined as a tuple $(s_i, R_i, a_i, s_{i+1}, o_i)$ that encapsulates a full state transition, where: s_i is the initial state before the action, represented by a UI screenshot. R_i is the agent’s cognitive reasoning trace, encapsulating its high-level strategic `Thought` and micro-level `Action Thought`. a_i is the specific `Action` executed by the agent in

Method	Level-1					Level-2					Level-3					Overall				
	SR		WPSR	MATCR	pATSR	SR		WPSR	MATCR	pATSR	SR		WPSR	MATCR	pATSR	SR		WPSR	MATCR	pATSR
	P@1	P@4				P@1	P@4				P@1	P@4				P@1	P@4			
PC-Agent																				
• Gemini-2.5-Pro	40.0	66.7	37.5	60.7	54.1	10.0	40.0	12.5	43.1	31.8	0.0	20.0	5.1	25.0	16.0	20.0	45.7	14.6	38.6	25.8
Mobile-Agent-e																				
• Gemini-2.5-Pro	46.7	100.0	57.1	75.0	70.0	10.0	<u>60.0</u>	20.0	52.5	41.9	0.0	<u>30.0</u>	9.3	24.6	17.2	22.9	<u>68.6</u>	23.2	44.5	31.6
LightManus																				
• Gemini-2.5-Pro	73.3	100.0	60.7	71.4	72.7	20.0	70.0	27.5	<u>53.1</u>	46.4	<u>10.0</u>	40.0	<u>16.9</u>	29.7	<u>23.9</u>	<u>40.0</u>	74.3	<u>29.9</u>	46.3	<u>37.1</u>
• Gemini-2.5-Flash	53.3	<u>93.3</u>	62.5	<u>74.1</u>	<u>70.7</u>	0.0	<u>60.0</u>	22.5	50.0	42.7	0.0	<u>30.0</u>	11.4	20.3	15.7	22.9	65.7	26.2	41.5	31.1
• GPT-o3	<u>60.0</u>	-	57.1	67.9	64.5	<u>30.0</u>	-	<u>30.0</u>	57.5	<u>49.6</u>	0.0	-	0.0	<u>32.2</u>	20.8	34.3	-	22.0	<u>48.0</u>	35.1
• Claude-sonnet-4	<u>60.0</u>	-	57.1	60.7	59.1	40.0	-	40.0	57.5	52.1	20.0	-	20.3	37.3	30.0	42.9	-	34.6	48.8	40.4
• Qwen3	40.0	60.0	28.6	50.0	44.5	0.0	30.0	10.0	32.5	25.2	<u>10.0</u>	10.0	5.1	16.1	10.2	20.0	37.1	11.8	28.7	19.2
• Qwen2.5-Max	<u>60.0</u>	86.7	50.0	58.9	55.9	10.0	20.0	10.0	30.6	23.2	<u>10.0</u>	20.0	4.7	11.0	7.5	31.4	48.6	16.3	27.8	18.6
• Qwen2.5-7B	6.7	26.7	7.1	18.8	15.7	10.0	10.0	2.5	21.9	15.1	0.0	10.0	2.1	14.8	7.8	5.7	17.1	3.3	17.9	11.0
• Qwen2.5-7B-RFT	33.3	73.3	28.6	40.2	37.0	10.0	20.0	7.5	23.8	17.3	0.0	20.0	4.7	16.5	10.0	17.1	42.8	10.8	24.0	15.8

Table 1: Main results on the NaturalGAIA, comparing the performance of different agent frameworks and foundation models across varying task difficulty levels. We evaluate using Success Rate (SR%) at Pass@1 and Pass@4, Weighted Pathway Success Rate (WPSR%), Mean Atomic Tasks Completion Ratio (MATCR%), and Positional-Weighted Atomic Tasks Success Rate (p-ATSR%). Our proposed LightManus framework consistently outperforms baseline agents. While Reinforcement Fine-Tuning (RFT) remarkably improves the smaller model’s performance, its effectiveness diminishes as task complexity increases.

that step. s_{i+1} is the new state reached after executing action a_i , also represented by a UI screenshot. o_i is the agent’s Observation and reflection on the outcome of the state transition from s_i to s_{i+1} .

This granular, structured format transparently captures the agent’s complete perception-reasoning-action loop, making it an ideal resource for process-supervision and model fine-tuning. We present detailed trajectory data samples and corresponding visualization execution trajectories in the Appendix E.

4.4 TRAJECTORY-BASED LLM ENHANCEMENT

The enhancement of the Qwen-2.5-VL-7B model was conducted via Trajectory-Reinforced Fine-Tuning (T-RFT), a methodological choice deliberately positioned between two common paradigms: Supervised Fine-Tuning (SFT) and end-to-end Reinforcement Learning (RL). While SFT offers a straightforward approach by mimicking expert demonstrations, it typically yields brittle policies that overfit to the specific, static trajectories seen during training. This reliance on behavioral cloning limits the agent’s ability to generalize to novel situations or recover from errors. Conversely, conventional end-to-end RL, though theoretically more robust, is often practically infeasible for complex GUI tasks. It necessitates extensive online exploration, entailing prohibitive computational overhead, significant communication latency, and long training durations.

T-RFT presents a pragmatic and resource-efficient alternative. By leveraging a high-quality, static dataset of human-verified trajectories, it circumvents the need for costly online interaction. Yet, unlike SFT, it moves beyond simple imitation by applying policy optimization principles to learn a more generalizable decision-making model. This allows the agent to distill robust navigational strategies from the expert data, rather than merely memorizing action sequences. Ultimately, this approach was chosen to specifically augment the model’s capacity to interpret and strategically act upon the sequential visual and textual inputs inherent to GUI environments.

5 EXPERIMENT

5.1 SETUP

Baselines. The performance of our proposed LightManus framework was benchmarked against two baseline agents: PC-Agent (PCA) for desktop environments and Mobile-Agent-e (MAe) for Android environments.

Foundation Models. Our evaluation includes a diverse range of state-of-the-art proprietary and open-source MLLMs. Proprietary models include Google’s Gemini series (Gemini-2.5-Pro, Gemini-2.5-Flash), OpenAI’s GPT-o3, and Anthropic’s Claude-sonnet-4. Open-source models

assessed are from the Qwen series, including Qwen3-235B-A22B, Qwen2.5-Max, and the base Qwen2.5-VL-7B model which was targeted for fine-tuning.

Training. The enhancement of the smaller model was achieved through T-RFT. Specifically, the Qwen2.5-VL-7B model was fine-tuned using our curated dataset comprising 278 human-verified trajectories. The training was executed on a high-performance computing cluster equipped with four NVIDIA A100 GPUs. Further training details can be found in Appendix F.

5.2 MAIN RESULTS

The experimental results, detailed in Table 1, reveal several key findings. Our proposed LightManus framework consistently demonstrates superior performance over baseline agents across all evaluated metrics. Among the foundation models, large-scale proprietary models, particularly Claude-sonnet-4 and Gemini-2.5-Pro, exhibit the most robust capabilities, achieving the highest success rates on complex, long-horizon tasks. While open-source models show promise, a discernible performance gap remains. Furthermore, T-RFT substantially enhances the capabilities of smaller models on basic tasks, yet this advantage diminishes markedly with increasing task complexity, highlighting the intrinsic limitations of model scale. The subsequent analysis dissects these findings by evaluating the distinct contributions of the agent framework, the foundation models, and the fine-tuning process.

Framework Efficacy: The Architectural Advantage of LightManus The LightManus framework demonstrates a significant advantage in GUI execution tasks. Driven by the same Gemini-2.5-Pro model, LightManus’s overall performance surpasses that of PCA and MAe. Specifically, its success rate (Pass@1) reaches 40.0%, which is considerably higher than PCA’s 20.0% and MAe’s 22.9%. Its WPSR is 29.9%, also outperforming PCA’s 14.6% and MAe’s 23.2%.

This superior performance stems from its higher execution reliability and long-horizon coherence. LightManus, powered by Gemini-2.5-Pro, shows substantial leads in both MATCR (46.3% vs. 39.2%/43.1%) and p-ATSR (37.1% vs. 25.8%/31.6%) metrics. A particularly compelling piece of evidence is that LightManus equipped with the smaller Gemini-2.5-Flash achieves a p-ATSR (31.1%) that is comparable to that of MAe equipped with the more powerful Gemini-2.5-Pro (31.6%). This result strongly indicates that the framework design of LightManus provides more robust execution support for complex, long-horizon tasks.

The Impact of LLM Scale on Performance The experimental results clearly reveal the critical impact of LLM scale on performance. Larger models, with their extensive knowledge bases, yield higher operational reliability. This is reflected in the MATCR metric: LightManus driven by Gemini-2.5-Pro (46.3%) outperforms its Gemini-2.5-Flash version (41.5%), indicating that larger models can execute atomic operations more stably.

However, a counter-intuitive phenomenon was observed in simple tasks: larger models occasionally exhibit a tendency to “over-plan”, breaking down simple instructions into excessively verbose steps. For instance, in Level-1 tasks, the agent equipped with Gemini-2.5-Flash slightly outperformed in WPSR (62.5% vs. 60.7%) and MATCR (74.1% vs. 71.4%). This suggests that for low-complexity tasks, the more concise planning strategy of smaller models can sometimes be more efficient.

Efficacy and Limitations of Small Model Specialization Fine-tuning smaller models highlights both their potential and inherent limitations. By fine-tuning Qwen2.5-VL-7B with our trajectory data, the model’s performance on basic tasks was significantly enhanced. In Level-1 tasks, its Pass@1 success rate increased nearly fivefold (6.7% \rightarrow 33.3%), WPSR nearly quadrupled (7.1% \rightarrow 28.6%), and MATCR also doubled (18.8% \rightarrow 40.2%).

However, this specialized capability gained through fine-tuning does not generalize well to complex scenarios. As task complexity increases, the model’s intrinsic capability ceiling becomes apparent. Its WPSR plummets from 28.6% at Level-1 to 7.5% at Level-2 and 4.7% at Level-3. This precipitous decline in performance indicates that core capabilities required for executing complex GUI tasks, such as long-range planning and memory, cannot be compensated for solely by fine-tuning on domain-specific data.

5.3 GENERALIZATION ACROSS APPLICATION SCENARIOS

Model	SR(%)
Qwen2.5-VL-7B	26.3
Qwen2.5-VL-7B-RFT	36.8

Table 2: Generalization performance on the A3 benchmark Chai et al. (2025). The evaluation is conducted on applications unseen during training to assess the model’s transferability.

As shown in table 2, to rigorously evaluate the model’s ability to generalize to new environments, we conducted a supplementary experiment on the A3 benchmark Chai et al. (2025). For this test, we selected six task subsets involving applications that are entirely disjoint from those present in our training trajectory data. This setup ensures a stringent assessment of transferable skills rather than memorized procedures. The results show that the base Qwen2.5-VL-7B model achieved a success rate of 26.3%. After T-RFT, the model’s performance on these novel applications rose to 36.8%, representing a substantial relative improvement of approximately 40%. This marked enhancement strongly suggests that the fine-tuning process did not lead to overfitting. Instead, it indicates the model has internalized more generalizable operational principles for GUI interaction, underscoring the value of our high-quality trajectory data in fostering robust and transferable agent capabilities.

5.4 ERROR ANALYSIS

As shown in Figure 2, we classify agent failures into three primary types: Planning & Reasoning Errors (PRE), which refer to when the agent formulates a logically flawed sequence of steps; Structural Compliance Errors (SCE), where the model fails to generate output in the required action format; and Execution Errors (EE), a category that encompasses the finer-grained Operational Errors (OE), Knowledge Deficits (KD), and Perceptual Errors (PE).

The error analysis reveals significant disparities in failure modes among the models. High-performance models, such as the Gemini series (-2.5-Pro and -2.5-Flash), exhibits low PRE and SCE rates, demonstrating robust capabilities in planning and structural compliance. Their failures are predominantly concentrated in EE, accounting for 47.8% and 54.3% of their respective outcomes. A deeper analysis of the EE category shows that the primary bottleneck for Gemini-2.5-Pro is OE (35.7%), highlighting a challenge in execution accuracy. For the smaller Gemini-2.5-Flash, its EE composition is more heavily skewed towards KD (19.3% vs. 5.0%), reflecting a trade-off between model scale and embedded knowledge.

In contrast, the Qwen2.5-VL-7B series of models exhibits a more fundamental bottleneck. The root cause of failure for the base model is a critical rate of SCE (46.4%), followed by EE (42.2%). While T-RFT substantially reduces the successful mitigation of KD (32.9% down to 11.4%). It is also noted that the SCE remains at a relatively high level of 47.9%. This indicates that while our fine-tuning method effectively imparts domain knowledge to improve execution, it has not yet fully resolved the inherent challenges in structural compliance found in certain smaller models.

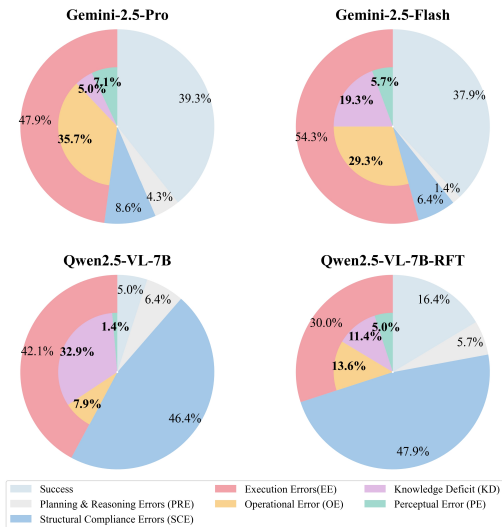


Figure 2: A comparative breakdown of failure modes for four vision-language models. The chart categorizes errors into three primary types: Planning & Reasoning Errors (PRE), Structural Compliance Errors (SCE), and Execution Errors (EE). The Execution Errors category is further subdivided into Operational Errors (OE), Knowledge Deficits (KD), and Perceptual Errors (PE).

6 CONCLUSION AND DISCUSSION

To address the critical challenges in evaluating GUI agents, we introduce NaturalGAIA, a novel benchmark based on “Causal Pathways”, a hierarchical LightManus architecture, and a high-quality trajectory dataset featuring diverse and self-correcting behaviors. The experiments demonstrate that even state-of-the-art models like Claude-sonnet-4 face formidable challenges on NaturalGAIA, revealing deficiencies in long-horizon planning. Meanwhile, although Trajectory-Reinforced Fine-Tuning (T-RFT) enhances the foundational performance of smaller models (e.g., Qwen2.5-VL-7B), their capability ceiling in complex tasks remains constrained: while fine-tuning supplements some world knowledge, it struggles to fundamentally improve their operational precision and perceptual capabilities, and has not yet fully addressed the underlying challenge in structural compliance. This set of limitations highlights the critical gap between domain-specific skills and generalizable reasoning abilities.

While NaturalGAIA provides a robust standard, it does not fully capture the dynamic nature of real-world interfaces. Future research can proceed in three directions: (1) developing agents with enhanced structural and logical reasoning abilities to overcome their innate bottlenecks; (2) creating dynamic benchmarks capable of auto-generating novel “Causal Pathways” to mitigate the risk of data contamination; (3) exploring advanced learning algorithms that can efficiently learn from self-correcting behaviors to build more resilient agents.

REFERENCES

- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, 2023. URL <https://arxiv.org/abs/2308.12966>.
- Yuxiang Chai, Hanhao Li, Jiayu Zhang, Liang Liu, Guangyi Liu, Guozhi Wang, Shuai Ren, Siyuan Huang, and Hongsheng Li. A3: Android agent arena for mobile gui agents, 2025. URL <https://arxiv.org/abs/2501.01149>.
- Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. SeeClick: Harnessing gui grounding for advanced visual gui agents, 2024. URL <https://arxiv.org/abs/2401.10935>.
- Thibault Le Sellier De Chezelles, Maxime Gasse, Alexandre Drouin, Massimo Caccia, Léo Boisvert, Megh Thakkar, Tom Marty, Rim Assouel, Sahar Omid Shayegan, Lawrence Keunho Jang, Xing Han Lù, Ori Yoran, Dehan Kong, Frank F. Xu, Siva Reddy, Quentin Cappart, Graham Neubig, Ruslan Salakhutdinov, Nicolas Chapados, and Alexandre Lacoste. The browsergym ecosystem for web agent research, 2025. URL <https://arxiv.org/abs/2412.05467>.
- DeepSeek-AI et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web, 2023. URL <https://arxiv.org/abs/2306.06070>.
- Zane Durante, Qiuyuan Huang, Naoki Wake, Ran Gong, Jae Sung Park, Bidipta Sarkar, Rohan Taori, Yusuke Noda, Demetri Terzopoulos, Yejin Choi, Katsushi Ikeuchi, Hoi Vo, Li Fei-Fei, and Jianfeng Gao. Agent ai: Surveying the horizons of multimodal interaction, 2024. URL <https://arxiv.org/abs/2401.03568>.
- Divyansh Garg, Shaun VanWeelden, Diego Caples, Andis Draguns, Nikil Ravi, Pranav Putta, Naman Garg, Tomas Abraham, Michael Lara, Federico Lopez, James Liu, Atharva Gundawar, Prannay Hebbar, Youngchul Joo, Jindong Gu, Charles London, Christian Schroeder de Witt, and Sumeet Motwani. Real: Benchmarking autonomous agents on deterministic simulations of real websites, 2025. URL <https://arxiv.org/abs/2504.11543>.
- Boyu Gou, Ruohan Wang, Boyuan Zheng, Yanan Xie, Cheng Chang, Yiheng Shu, Huan Sun, and Yu Su. Navigating the digital world as humans do: Universal visual grounding for gui agents, 2025. URL <https://arxiv.org/abs/2410.05243>.

- 540 Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan,
541 and Dong Yu. Webvoyager: Building an end-to-end web agent with large multimodal models,
542 2024. URL <https://arxiv.org/abs/2401.13919>.
- 543
544 Yanheng He, Jiahe Jin, and Pengfei Liu. Efficient agent training for computer use, 2025. URL
545 <https://arxiv.org/abs/2505.13909>.
- 546 Wenyi Hong, Weihao Wang, Qingsong Lv, Jiazhen Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan
547 Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent:
548 A visual language model for gui agents, 2024. URL [https://arxiv.org/abs/2312.](https://arxiv.org/abs/2312.08914)
549 08914.
- 550
551 Zhen Huang, Haoyang Zou, Xuefeng Li, Yixiu Liu, Yuxiang Zheng, Ethan Chern, Shijie Xia, Yiwei
552 Qin, Weizhe Yuan, and Pengfei Liu. O1 replication journey – part 2: Surpassing o1-preview
553 through simple distillation, big progress or bitter lesson?, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2411.16489)
554 [abs/2411.16489](https://arxiv.org/abs/2411.16489).
- 555 Wenjia Jiang, Yangyang Zhuang, Chenxi Song, Xu Yang, and Chi Zhang. Appagentx: Evolving
556 gui agents as proficient smartphone users, 2025. URL [https://arxiv.org/abs/2503.](https://arxiv.org/abs/2503.02268)
557 02268.
- 558
559 Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Lim, Po-Yu Huang, Graham
560 Neubig, Shuyan Zhou, Russ Salakhutdinov, and Daniel Fried. VisualWebArena: Evaluating
561 multimodal agents on realistic visual web tasks. In Lun-Wei Ku, Andre Martins, and Vivek
562 Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational*
563 *Linguistics (Volume 1: Long Papers)*, pp. 881–905, Bangkok, Thailand, August 2024a. As-
564 sociation for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.50. URL [https:](https://aclanthology.org/2024.acl-long.50/)
565 [//aclanthology.org/2024.acl-long.50/](https://aclanthology.org/2024.acl-long.50/).
- 566
567 Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham
568 Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating
569 multimodal agents on realistic visual web tasks, 2024b. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2401.13649)
570 2401.13649.
- 571
572 Kaixin Li, Ziyang Meng, Hongzhan Lin, Ziyang Luo, Yuchen Tian, Jing Ma, Zhiyong Huang, and
573 Tat-Seng Chua. Screenspot-pro: Gui grounding for professional high-resolution computer use,
574 2025a. URL <https://arxiv.org/abs/2504.07981>.
- 575
576 Zhangheng Li, Keen You, Haotian Zhang, Di Feng, Harsh Agrawal, Xiujun Li, Mohana
577 Prasad Sathya Moorthy, Jeff Nichols, Yinfei Yang, and Zhe Gan. Ferret-ui 2: Mastering uni-
578 versal user interface understanding across platforms, 2025b. URL [https://arxiv.org/](https://arxiv.org/abs/2410.18967)
579 [abs/2410.18967](https://arxiv.org/abs/2410.18967).
- 580
581 Guangyi Liu, Pengxiang Zhao, Liang Liu, Yaxuan Guo, Han Xiao, Weifeng Lin, Yuxiang Chai, Yue
582 Han, Shuai Ren, Hao Wang, Xiaoyu Liang, Wenhao Wang, Tianze Wu, Linghao Li, Hao Wang,
583 Guanqing Xiong, Yong Liu, and Hongsheng Li. Llm-powered gui agents in phone automation:
584 Surveying progress and prospects, 2025a. URL <https://arxiv.org/abs/2504.19838>.
- 585
586 Haowei Liu, Xi Zhang, Haiyang Xu, Yuyang Wanyan, Junyang Wang, Ming Yan, Ji Zhang, Chun-
587 feng Yuan, Changsheng Xu, Weiming Hu, and Fei Huang. Pc-agent: A hierarchical multi-
588 agent collaboration framework for complex task automation on pc, 2025b. URL [https:](https://arxiv.org/abs/2502.14282)
589 [//arxiv.org/abs/2502.14282](https://arxiv.org/abs/2502.14282).
- 590
591 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,
592 Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui
593 Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie
594 Tang. Agentbench: Evaluating llms as agents, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2308.03688)
595 2308.03688.
- 596
597 Grégoire Mialon, Clémentine Fourier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas
598 Scialom. Gaia: a benchmark for general ai assistants, 2023. URL [https://arxiv.org/](https://arxiv.org/abs/2311.12983)
599 [abs/2311.12983](https://arxiv.org/abs/2311.12983).

- 594 Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke
595 Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. *s1: Simple test-time*
596 *scaling*, 2025. URL <https://arxiv.org/abs/2501.19393>.
- 597
598 OpenAI et al. *Gpt-4 technical report*, 2024a. URL <https://arxiv.org/abs/2303.08774>.
- 599
600 OpenAI et al. *Openai o1 system card*, 2024b. URL <https://arxiv.org/abs/2412.16720>.
- 601
602 Jiahao Qiu, Xuan Qi, Tongcheng Zhang, Xinzhe Juan, Jiacheng Guo, Yifu Lu, Yimin Wang,
603 Zixin Yao, Qihan Ren, Xun Jiang, Xing Zhou, Dongrui Liu, Ling Yang, Yue Wu, Kaixuan
604 Huang, Shilong Liu, Hongru Wang, and Mengdi Wang. *Alita: Generalist agent enabling scal-*
605 *able agentic reasoning with minimal predefinition and maximal self-evolution*, 2025. URL
606 <https://arxiv.org/abs/2505.20286>.
- 607
608 Qwen et al. *Qwen2.5 technical report*, 2025. URL <https://arxiv.org/abs/2412.15115>.
- 609
610 Hanoona Rasheed, Muhammad Maaz, Sahal Shaji Mullappilly, Abdelrahman Shaker, Salman Khan,
611 Hisham Cholakkal, Rao M. Anwer, Erix Xing, Ming-Hsuan Yang, and Fahad S. Khan. *Glamm:*
612 *Pixel grounding large multimodal model*, 2024. URL <https://arxiv.org/abs/2311.03356>.
- 613
614 Christopher Rawles, Sarah Clinckemaiellie, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Marybeth
615 Fair, Alice Li, William Bishop, Wei Li, Folawiyo Campbell-Ajala, Daniel Toyama, Robert Berry,
616 Divya Tyamagundlu, Timothy Lillicrap, and Oriana Riva. *Androidworld: A dynamic benchmark-*
617 *ing environment for autonomous agents*, 2024.
- 618
619 Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and
620 Shunyu Yao. *Reflexion: Language agents with verbal reinforcement learning*, 2023. URL
621 <https://arxiv.org/abs/2303.11366>.
- 622
623 Fei Tang, Haolei Xu, Hang Zhang, Siqi Chen, Xingyu Wu, Yongliang Shen, Wenqi Zhang, Guiyang
624 Hou, Zeqi Tan, Yuchen Yan, Kaitao Song, Jian Shao, Weiming Lu, Jun Xiao, and Yueting
625 Zhuang. *A survey on (m)llm-based gui agents*, 2025. URL <https://arxiv.org/abs/2504.13865>.
- 626
627 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan,
628 and Anima Anandkumar. *Voyager: An open-ended embodied agent with large language models*,
629 2023. URL <https://arxiv.org/abs/2305.16291>.
- 630
631 Junyang Wang, Haiyang Xu, Haitao Jia, Xi Zhang, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang,
632 and Jitao Sang. *Mobile-agent-v2: Mobile device operation assistant with effective navigation via*
633 *multi-agent collaboration*, 2024. URL <https://arxiv.org/abs/2406.01014>.
- 634
635 Zhenhailong Wang, Haiyang Xu, Junyang Wang, Xi Zhang, Ming Yan, Ji Zhang, Fei Huang, and
636 Heng Ji. *Mobile-agent-e: Self-evolving mobile assistant for complex tasks*. *arXiv preprint*
637 *arXiv:2501.11733*, 2025.
- 638
639 Yifan Xu, Xiao Liu, Xueqiao Sun, Siyi Cheng, Hao Yu, Hanyu Lai, Shudan Zhang, Dan Zhang,
640 Jie Tang, and Yuxiao Dong. *Androidlab: Training and systematic benchmarking of android au-*
641 *tonomous agents*, 2024. URL <https://arxiv.org/abs/2410.24024>.
- 642
643 Sherry Yang, Ofir Nachum, Yilun Du, Jason Wei, Pieter Abbeel, and Dale Schuurmans. *Foundation*
644 *models for decision making: Problems, methods, and opportunities*, 2023. URL [https://](https://arxiv.org/abs/2303.04129)
645 arxiv.org/abs/2303.04129.
- 646
647 Suyu Ye, Haojun Shi, Darren Shih, Hyokun Yun, Tanya Roosta, and Tianmin Shu. *Realwebassist:*
648 *A benchmark for long-horizon web assistance with real-world users*, 2025a. URL [https://](https://arxiv.org/abs/2504.10445)
649 arxiv.org/abs/2504.10445.
- 650
651 Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie Xia, and Pengfei Liu. *Limo: Less is more*
652 *for reasoning*, 2025b. URL <https://arxiv.org/abs/2502.03387>.

648 Chaoyun Zhang, He Huang, Chiming Ni, Jian Mu, Si Qin, Shilin He, Lu Wang, Fangkai Yang,
649 Pu Zhao, Chao Du, Liqun Li, Yu Kang, Zhao Jiang, Suzhen Zheng, Rujia Wang, Jiaxu Qian,
650 Minghua Ma, Jian-Guang Lou, Qingwei Lin, Saravan Rajmohan, and Dongmei Zhang. Ufo2:
651 The desktop agentos, 2025. URL <https://arxiv.org/abs/2504.14603>.

652 Wenqi Zhang, Ke Tang, Hai Wu, Mengna Wang, Yongliang Shen, Guiyang Hou, Zeqi Tan, Peng Li,
653 Yueting Zhuang, and Weiming Lu. Agent-pro: Learning to evolve via policy-level reflection and
654 optimization, 2024. URL <https://arxiv.org/abs/2402.17574>.

656 Xinlu Zhang, Yujie Lu, Weizhi Wang, An Yan, Jun Yan, Lianke Qin, Heng Wang, Xifeng Yan,
657 William Yang Wang, and Linda Ruth Petzold. Gpt-4v(ision) as a generalist evaluator for vision-
658 language tasks, 2023. URL <https://arxiv.org/abs/2311.01361>.

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

664 A ETHICS STATEMENT

666 This work adheres to the ICLR Code of Ethics. In this study, no human subjects or animal ex-
667 perimentation was involved. All datasets used, were sourced in compliance with relevant usage
668 guidelines, ensuring no violation of privacy. We have taken care to avoid any biases or discrimi-
669 natory outcomes in our research process. No personally identifiable information was used, and no
670 experiments were conducted that could raise privacy or security concerns. We are committed to
671 maintaining transparency and integrity throughout the research process.

673 B REPRODUCIBILITY STATEMENT

675 We have made every effort to ensure that the results presented in this paper are reproducible. All
676 code and datasets have been made publicly available in an anonymous repository to facilitate repli-
677 cation and verification. The experimental setup, including training steps, model configurations, and
678 hardware details, is described in detail in the paper.

679 We believe these measures will enable other researchers to reproduce our work and further advance
680 the field.

682 C LLM USAGE

684 Large Language Models (LLMs) were used to aid in the writing and polishing of the manuscript.
685 Specifically, we used an LLM to assist in refining the language, improving readability, and ensuring
686 clarity in various sections of the paper. The model helped with tasks such as sentence rephrasing,
687 grammar checking, and enhancing the overall flow of the text.

689 It is important to note that the LLM was not involved in the ideation, research methodology, or
690 experimental design. All research concepts, ideas, and analyses were developed and conducted by
691 the authors. The contributions of the LLM were solely focused on improving the linguistic quality
692 of the paper, with no involvement in the scientific content or data analysis.

693 The authors take full responsibility for the content of the manuscript, including any text generated
694 or polished by the LLM. We have ensured that the LLM-generated text adheres to ethical guidelines
695 and does not contribute to plagiarism or scientific misconduct.

697 D DETAILS OF NATURALGAIA

699 D.1 DESIGN PRINCIPLES OF NATURALGAIA (LEVEL-3)

701 Taking the Level-3 build as an example, the design guidelines we provide to Task builders are as
follows:

702 **Part I: Task Structure and Tool Application**
703

- 704
- 705 • **Number and Identification of Atomic Tasks:** Each main task must be designed to contain
706 six to seven clearly distinguishable atomic tasks. Each atomic task should be clearly labeled
707 and distinguished in the task description document or related schematic diagram.
 - 708 • **Multi-tool Collaboration Requirements:** Ensure that at least three different mobile ap-
709 plications (APPs) are called to complete each main task. In terms of application selection,
710 applications on the Android platform should be selected.

711 **Part II: Core Characteristics of Atomic Task Results**
712

- 713
- 714 • **Textual Presentation of Results:** The final results of the constructed atomic tasks must be
715 clearly presented in text form to support subsequent reliability assessments.
 - 716 • **Uniqueness and Certainty:** The output results generated after the execution of each
717 atomic task must be unique and clearly defined, and can be verified objectively and un-
718 ambiguously, and there should not be multiple possible correct answers.
 - 719 • **Exclusion of Timeliness Effects:** The results of atomic tasks should be designed to be
720 time-invariant, aiming to avoid changes, failures, or deviations from expected standards
721 due to the passage of time.
 - 722 • **Processing procedures for time-sensitive information:** If the task scenario inevitably
723 involves critical information with timeliness (for example, real-time stock prices, limited-
724 time promotion information, dynamically updated news lists, etc.), the information must
725 be solidified the moment it is captured by the atomic task. This processing requires that the
726 information content be fully recorded in a local static file (e.g.: screenshots saved as PNG,
727 text information saved as TXT/PDF, table data saved as XLSX, complex documents saved
728 as DOCX, etc.). At the same time, the atomic task must provide specific time constraints
729 (e.g.: query the price of Apple16 on Apple’s official website on April 11, 2025). The local
730 static file constitutes the deterministic and verifiable output result of this atomic task.

731 **Part III: Logical dependencies between atomic tasks**
732

- 733
- 734 • **Mandatory information flow dependencies:** Strict sequential dependencies must be es-
735 tablished between the designed atomic tasks. Specifically, the output result z of the previous
736 atomic task must be used as the core input information required for the execution of the next
737 atomic task.
 - 738 • **Failed step-by-step interruption mechanism:** If the previous atomic task fails to suc-
739 cessfully obtain its expected correct result, all subsequent atomic tasks that rely on this
740 result should not be able to continue or be completed correctly. The task process must be
741 designed based on this standard.

742 **Part IV: Clarification of application tools and ground truth**
743

- 744
- 745 • **Application of atomic tasks:** In the detailed description of each atomic task, the name of
746 the specific mobile APP expected to be used to complete the task must be clearly specified.
 - 747 • **Detailed description of the path to obtain the ground truth:** The specific operation path
748 or method used by each atomic task constructed to obtain its ground truth must be described
749 in detail to facilitate verification and confirmation by the verifier.

750 **Part V: Selection criteria for applications**
751

- 752
- 753 • **Commonality priority principle:** When selecting applications for tasks, be sure to give
754 priority to common applications with high user penetration and widespread use in the app
755 store¹² (as of April 2025, ranked within the top 50).

Table 3: An example of a task definition from Stage 1. The table specifies the high-level task, its difficulty level, the step-by-step Causal Pathway for human execution, the corresponding ground truth answers, and the tools involved.

Task	Search for Leonardo Dicaprio on Douban and tell me what movie the actor is starring in in 2023? Use IMDB to look up the movie’s profile and tell me what era the story takes place in?
Level	1
Steps(human)	<ol style="list-style-type: none"> 1. Open the Douban app. 2. Type “Leonardo Dicaprio” in the search box and click Search. 3. In the search results, click on Leonardo Dicaprio’s actor entry to enter his personal homepage. 4. View his film and television works list and find the works released in “2023”. 5. Record the movie name of the 2023 work. (Ground Truth 1: Killers of the Flower Moon) 6. Open the IMDB app. 7. Type the movie name recorded from the Douban app (Killers of the Flower Moon) in the search box and click Search. 8. Click the corresponding movie entry in the search results to enter the movie details page. 9. Scroll down to the movie profile section. 10. Read, find and record the era information mentioned in the profile (Ground Truth 2: 1920s).
Ground Truth	<ol style="list-style-type: none"> 1. Killers of the Flower Moon 2. 1920s
Tools	Douban app IMDB APP
Number of Atomic Tasks	2

D.2 DETAILS OF TASK CREATION

To ensure the quality, feasibility, and determinism of each task, a rigorous three-stage process was designed and implemented for the construction of NaturalGAIA, executed by a team of four annotators. This multi-stage protocol ensures the robustness of the tasks and the reliability of their outcomes. The process is detailed as follows:

Stage 1: Task Design and Causal Pathway Definition. This initial phase was led by Annotator 1, who designed the Causal Pathway and its constituent atomic tasks. The design was governed by core principles mandating that the output (i.e., the answer) for each task node be unique, deterministic, and temporally stable. This foundation ensures the long-term validity of the benchmark. Annotator 1 also developed reference procedures to confirm the pathway’s initial feasibility on the Android platform. An example of a complete task definition produced in this stage, detailing the task, its constituent steps (the Causal Pathway), and ground truth, is presented in Table 3.

Stage 2: Independent Verification and Cross-Platform Screening. Annotators 2 and 3 served as independent validators. They first followed the reference procedures to traverse the designed Causal Pathway on Android. Subsequently, they verified the pathway’s applicability on PC (Windows) platforms. This stage rigorously evaluated the practical feasibility and output determinism of the tasks, ensuring cross-platform reproducibility. The verification results for the example task, confirming answer consistency across Android and Windows platforms, are shown in Table 4.

¹<https://play.google.com/store>.

²<https://m.app.mi.com/>.

Table 4: Cross-platform verification for the example task, corresponding to Stage 2. This table demonstrates the task’s reproducibility by confirming that identical ground truth answers are obtained on both Android and Windows platforms.

Operating System		Android	Windows
Task		Search for Leonardo Dicaprio on Douban and tell me what movie the actor is starring in in 2023? Use IMDB to look up the movie’s profile and tell me what era the story takes place in?	
Level		1	
Ground Truth		1. Killers of the Flower Moon 2. 1920s	
Tools		Douban app IMDB APP	https://www.douban.com/ https://www.imdb.com/
Answer Matching	Atomic Task 1	✓	✓
	Atomic Task 2	✓	✓

Stage 3: Blind Feasibility Assessment and Human Baseline Establishment. This phase was conducted by Annotator 4 to assess task clarity and establish a human performance baseline. To ensure an unbiased evaluation, Annotator 4 was provided only with the high-level task description, without access to the reference solution or the intended Causal Pathway. This setup simulates the process of an agent discovering and navigating the pathway based solely on instructions, thereby providing a crucial human baseline for evaluating agent autonomy.

D.3 DETAILS OF DATA STATISTICS

Table 5 presents the detailed statistical characteristics of the NaturalGAIA benchmark at different difficulty levels. The table lists the total number of entries (Num) in each difficulty interval from level 1 to level 5 and the Overall case. For each difficulty level, we further quantified the average number (Avg) and distribution Range of Tools used by annotators and Atomic tasks) performed during the benchmark construction process. In NaturalGAIA, there is a positive correlation between the difficulty level of the task and the average number of tools and atomic tasks required, which indicates that tasks with higher difficulty usually correspond to more complex processing flows and operation sequences.

Table 5: Statistics of NaturalGAIA

Difficulty level	Num	Tools		Atomic tasks	
		Avg	Range	Avg	Range
1	133	1.5	1~2	1.9	1~2
2	92	2.3	2~4	3.7	3~4
3	51	4.4	3~7	5.7	3~7
Overall	276	2.6	1~7	3.1	1~7

D.4 TASK EXAMPLE OF NATURALGAIA

Figure 3 presents concrete examples of tasks corresponding to the three difficulty levels. These examples illustrate the increasing requirements in terms of the number of atomic tasks, the diversity of applications involved, and the complexity of planning and execution across the levels.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

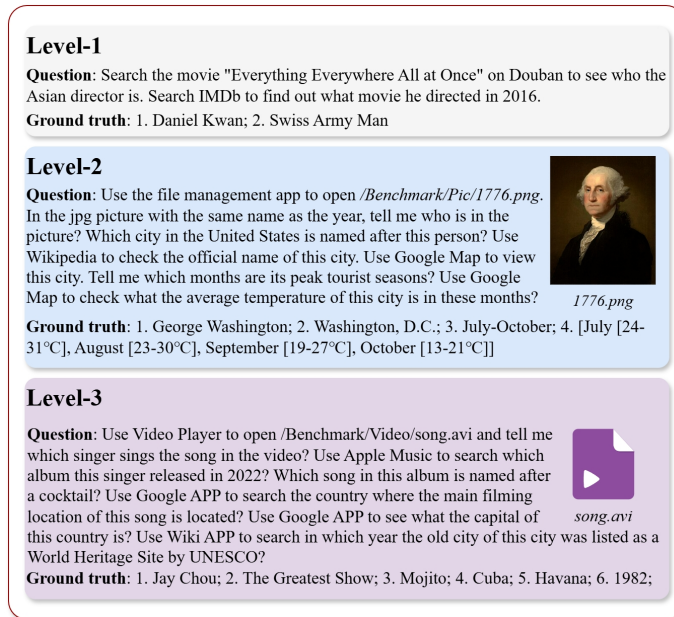


Figure 3: Task examples of different difficulty levels in NaturalGAIA

E DETAILS OF EXECUTION TRAJECTORY

Each data entry in our collection is a JSON object designed to provide a comprehensive and transparent record of an agent’s problem-solving process for a given task. The structure is organized with the following top-level keys:

- **messages:** This key holds an array that documents the agent’s interaction flow. It begins with the initial, high-level task instruction (e.g., “Using Wikipedia, search for Jay Chou and tell me which album he released in 2000?”), followed by the agent’s detailed, step-by-step execution trace.
- **solution:** This key contains a string with the ground-truth answer for the task, encapsulated within an `<answer>` tag (e.g., `<answer>Jay</answer>`). This allows for direct and automated evaluation of the task outcome.
- **images:** This key holds an array of strings, where each string is a file path to a screenshot captured during the interaction. These images provide the visual context for each step of the agent’s trajectory.

The core of our dataset lies within the **execution trace**, which documents the agent’s perception-reasoning-action loop for each step. A single step is structured as follows:

- **think Block:** This block encapsulates the agent’s cognitive process before taking an action. It is crucial for understanding the agent’s intent and planning. It contains:
 - **<image> Placeholder:** Marks where the corresponding pre-action screenshot is presented to the model.
 - **Thought:** A paragraph of text detailing the agent’s high-level strategic thinking. It assesses the overall progress, recalls the main objective, and formulates a sub-goal for the current step.
 - **Action Thought:** A paragraph focused on low-level, tactical reasoning. Grounded in the visual information from the `<image>`, it identifies the specific UI element to interact with (e.g., a button or text field) and determines the **type of action** to perform on it (e.g., tap or type text).

- 918
- 919 • **Action Block:** This block contains a concise, human-readable natural language description of the action the agent decided to execute (e.g., “Tap on ‘Search Wikipedia’ to activate the keyboard.”).
 - 920
 - 921 • **Observation Block:** This block documents the agent’s reflection after the action has been executed. It includes:
 - 922
 - 923 – **<image> Placeholder:** Marks where the post-action screenshot is presented, showing the result of the action.
 - 924
 - 925 – **Observation:** A textual summary where the agent confirms the outcome of its action and updates its understanding of the task status, setting the stage for the next step.
 - 926
 - 927

928 This detailed, multi-faceted structure is explicitly designed to provide rich, process-level supervision signals, enabling models to learn not just what to do, but also why and how to do it based on visual evidence and logical reasoning. The trajectory data sample and visualization process are shown in the listings 1 to 3 and figs. 4 to 6 respectively.

```

932 {
933   "images": [
934     "image_17-44-23_1_1.jpg",
935     "image_17-44-23_1_2.jpg",
936     "image_17-44-23_1_3.jpg",
937     "image_17-44-23_1_4.jpg",
938     "image_17-44-23_1_5.jpg"
939   ],
940   "messages": [
941     {
942       "role": "user",
943       "content": "Use Wikipedia to search for Jay Chou, and
944         tell me which album he released in 2000?"
945     },
946     {
947       "role": "assistant",
948       "content": "=== Step1 ==="
949     }
950   ]
951 }

```

949 <think>

950 <image_17-44-23_1_1.jpg>

951

952 ****Thought****

953 The user wants to find the name of Jay Chou’s album released in

954 2000 on Wikipedia. I need to open the Wikipedia app, search

955 for “Jay Chou”, navigate to his discography or album section,

956 and then find the album released in 2000.

957

958 ****Action Thought****

959 The current subgoal is to open the Wikipedia app. Looking at the

960 screen information, I can see “Wikipedia” at coordinates [664,

961 1023]. This is the app I need to open.

962

963 ****Action****

964 Open the Wikipedia app to start searching for “Jay Chou”.

965

966 </think>

967 <Observation>

968

969 <image17-44-23_1_2.jpg >

970

971 ****Observation****

1. Open the Wikipedia app.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025



Figure 4: Image from image_17-44-23_1.1 . jpg to image_17-44-23_1.5 . jpg

</Observation >

=== Step2 ===

1026
1027 <think>
1028
1029 <image17-44-23_1_2.jpg >
1030
1031 ****Thought****
1032 The Wikipedia app is open. The next step is to search for "Jay
1033 Chou". I can see a search bar labeled "Search Wikipedia" at
1034 the top of the screen. I will tap on this search bar to
1035 initiate the search.

1036 ****Action Thought****
1037 The current subgoal is to search for "Jay Chou". Looking at the
1038 screenshot, I see a search bar with the text "Search Wikipedia
1039 ". I need to tap on this search bar to activate the keyboard
1040 and then type "Jay Chou". The coordinates for "Search
1041 Wikipedia" are [323, 387].

1042 ****Action****
1043 Tap on the "Search Wikipedia" bar to activate the keyboard and
1044 prepare for typing the search query.
1045
1046 </think>
1047
1048 <Observation>
1049
1050 <image_17-44-23_1_3.jpg >
1051
1052 ****Observation****
1053 1. Open the Wikipedia app. 2. Search for "Jay Chou".
1054
1055 </Observation>
1056
1057 ===Step3===
1058
1059 <think>
1060
1061 <image_17-44-23_1_3.jpg >
1062
1063 ****Thought****
1064 The current screen shows the Wikipedia search interface, and "Jay
1065 Chou" is listed under "Recent searches". This means the search
1066 term has been entered, and I can proceed by tapping on "Jay
1067 Chou" to view the search results. This completes the "Search
1068 for 'Jay Chou'" subgoal. The next step is to navigate to the
1069 album/discography section.

1070 ****Action Thought****
1071 The current screen shows the Wikipedia search interface. The user's
1072 instruction is to search for "Jay Chou". I have already
1073 opened the Wikipedia app and tapped on the search bar. Now I
1074 need to type "Jay Chou" into the search bar. Looking at the
1075 screen information, "Jay Chou" is listed under "Recent
1076 searches". This means I don't need to type it, I can just tap
1077 on the existing "Jay Chou" entry.

1078 ****Action****
1079 Tap on "Jay Chou" from the recent searches to perform the search.

```

1080 </think>
1081
1082 <Observation>
1083
1084 <image_17-44-23_1_4.jpg >
1085
1086 **Observation**
1087 1. Open the Wikipedia app. 2. Search for "Jay Chou". 3. Navigate
1088     to the album/discography section of the search result.
1089
1090 </Observation>
1091
1092 ===Step4===
1093
1094 <think>
1095
1096 <image_17-44-23_1_4.jpg >
1097
1098 **Thought**
1099 The current screen shows the search results for "Jay Chou". The
1100     goal is to find the album released in 2000. To do this, I need
1101     to navigate to the album/discography section. Looking at the
1102     options, "Jay Chou Discography" seems to be the most relevant
1103     link to find a list of his albums. "Jay (album)" is also an
1104     album, but it's a specific one, and I need to find the one
1105     released in 2000, so a discography list would be more
1106     comprehensive. Therefore, clicking on "Jay Chou Discography"
1107     is the next logical step.
1108
1109 **Action Thought**
1110 The user wants to find the album released by Jay Chou in 2000. The
1111     current screen shows search results related to "Jay Chou".
1112     Among the results, "Jay Chou Discography" seems to be the most
1113     relevant link to find album information. I will tap on this
1114     link to navigate to the discography section.
1115
1116 **Action**
1117 Tap on "Jay Chou Discography" to navigate to the discography
1118     section, which should contain information about his albums and
1119     their release dates.
1120
1121 </think>
1122
1123 <Observation>
1124
1125 <image_17-44-23_1_5.jpg >
1126
1127 **Observation**
1128 1. Open the Wikipedia app. 2. Search for "Jay Chou". 3. Navigate
1129     to the album/discography section of the search result. 4.
1130     Identify the album released in 2000.
1131
1132 </Observation>"
1133     }
1134     ],
1135     "solution": "<answer> Jay </answer>"
1136 }

```

Listing 1: Example of Execution Trajectory 1

1134
1135
1136
1137
1138
1139
1140
1141
1142
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

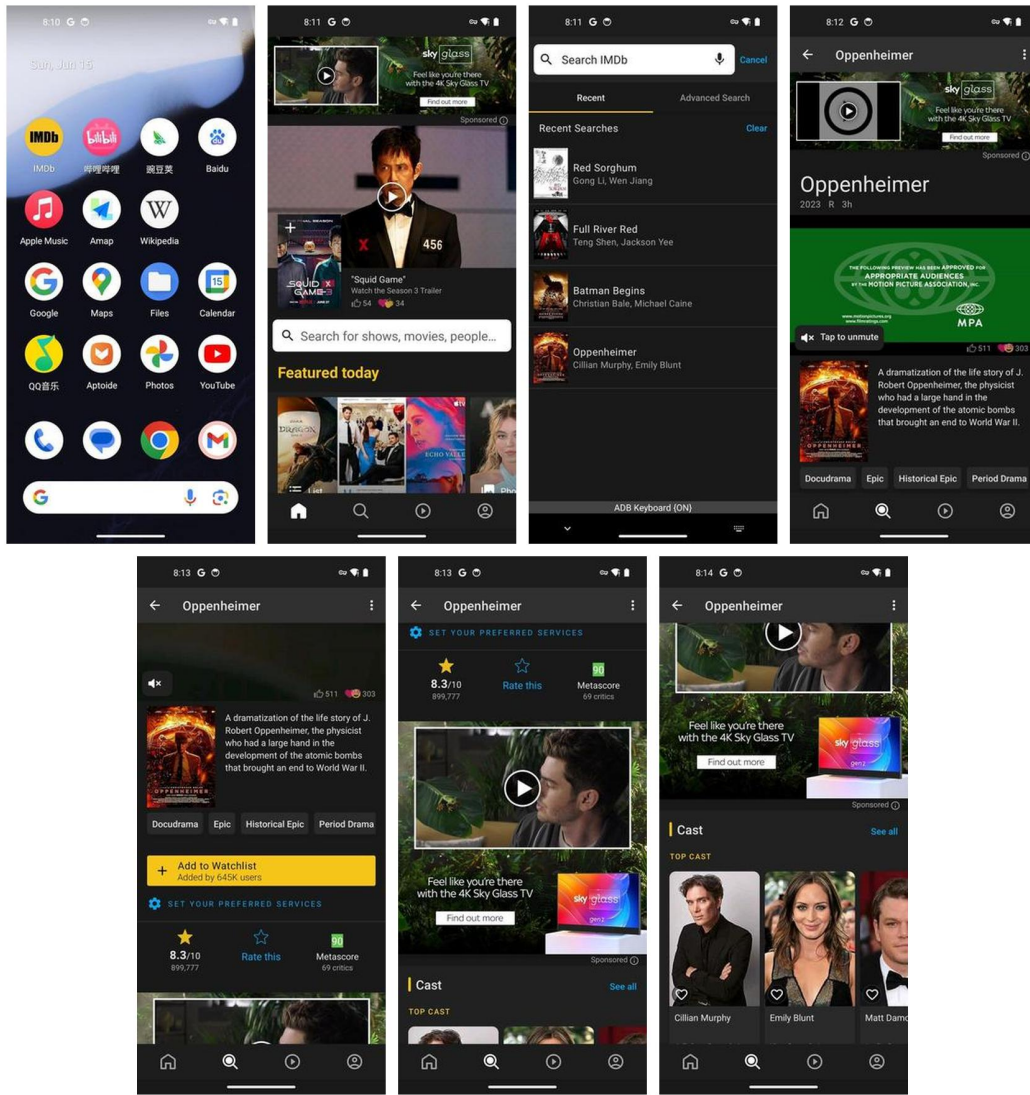


Figure 5: Image from image_10-41.1.1 .jpg to image_10-41.1.7 .jpg

```
{
  "images": [
    "image_-10-41_1.1 .jpg",
    "image_-10-41_1.2 .jpg",
    "image_-10-41_1.3 .jpg",
    "image_-10-41_1.4 .jpg",
    "image_-10-41_1.5 .jpg",
    "image_-10-41_1.6 .jpg",
    "image_-10-41_1.7 .jpg"
  ],
  "messages": [
    {
      "role": "user",
      "content": "Please use the IMDb app to search for 'Oppenheimer' and tell me who its first lead actor is?"
    }
  ]
}
```

1188 "role": "assistant",
 1189 "content": "===Step1===
 1190
 1191 <think>
 1192
 1193 <image_-10-41_1_1.jpg>
 1194
 1195 ****Thought****
 1196 The user wants to search for the movie "Oppenheimer" in the IMDb
 1197 app and find its first lead actor. This task can be broken
 1198 down into the following steps: 1. Open the IMDb app. 2. Search
 1199 for "Oppenheimer". 3. Navigate to the movie's details page.
 1200 4. Identify the first lead actor.
 1201
 1202 ****Action Thought****
 1203 The user wants to search for a movie in the IMDb app. The first
 1204 step is to open the IMDb app. Looking at the screen
 1205 information, I can see an icon labeled "IMDb" at coordinates
 1206 [162, 570]. This is the app I need to open.
 1207
 1208 ****Action****
 1209 Open the IMDb app to begin the search for the movie "Oppenheimer".
 1210
 1211 </think>
 1212
 1213 <Observation>
 1214
 1215 <image_-10-41_1_2.jpg>
 1216
 1217 ****Observation****
 1218 1. Open the IMDb app. 2. Search for the movie "Oppenheimer".
 1219
 1220 </Observation>
 1221
 1222 **===Step2===**
 1223
 1224 <think>
 1225
 1226 <image_-10-41_1_2.jpg>
 1227
 1228 ****Thought****
 1229 The current screen is the IMDb app's home screen. The previous
 1230 subgoal was to open the IMDb app, which is already completed.
 1231 The next step in the plan is to search for "Oppenheimer"
 1232 within the IMDb app. I can see a search bar labeled "Search
 1233 for shows, movies, people..." in the middle of the screen. I
 1234 should tap on this search bar to initiate the search.
 1235
 1236 ****Action Thought****
 1237 The current subgoal is to search for "Oppenheimer" within the IMDb
 1238 app. The screenshot shows the IMDb app's home screen, and
 1239 there is a search bar labeled "Search for shows, movies,
 1240 people..." in the middle of the screen. I need to tap on this
 1241 search bar to activate the keyboard and then type "Oppenheimer".
 1242
 1243 ****Action****
 1244 Tap on the search bar to activate the keyboard and prepare for
 1245 typing "Oppenheimer".

1242
1243 </think>
1244
1245 <Observation>
1246
1247 <image_-10-41_1_3.jpg>
1248
1249 **** Observation****
1250 1. Open the IMDb app. 2. Search for the movie "Oppenheimer". 3.
1251 Type "Oppenheimer" into the search bar.
1252 </Observation>
1253
1254 ===Step3===
1255
1256 <think>
1257
1258 <image_-10-41_1_3.jpg>
1259
1260 **** Thought****
1261 The current screen shows the search results within the IMDb app.
1262 The movie "Oppenheimer" is already listed under "Recent
1263 Searches". Since the goal is to search for "Oppenheimer" and
1264 then find its first lead actor, the next logical step is to
1265 tap on the "Oppenheimer" search result to view its details.
1266 This will complete the "Search for 'Oppenheimer' within the
1267 IMDb app" subgoal and move towards the next one.
1268 **** Action Thought****
1269 The current subgoal is to tap on the search result for "
1270 Oppenheimer" to view its details. Looking at the screen
1271 information, "Oppenheimer" is listed as a recent search result
1272 at coordinates [331, 1508]. I should tap on this to proceed.
1273
1274 **** Action****
1275 Tap on the "Oppenheimer" search result to view its details page.
1276 </think>
1277
1278 <Observation>
1279
1280 <image_-10-41_1_4.jpg>
1281
1282 **** Observation****
1283 1. Open the IMDb app. 2. Search for the movie "Oppenheimer". 3.
1284 Type "Oppenheimer" into the search bar. 4. Tap on the search
1285 result for "Oppenheimer" to view its details.
1286 </Observation>
1287
1288 ===Step4===
1289
1290 <think>
1291
1292 <image_-10-41_1_4.jpg>
1293
1294 **** Thought****
1295

1296 The current screen shows the movie details page for "Oppenheimer".
 1297 The goal is to find the first actor listed. I need to scroll
 1298 down to find the cast section.
 1299

1300 ****Action Thought****
 1301 The current screen shows the movie details page for "Oppenheimer".
 1302 The goal is to identify the first actor listed. I need to
 1303 scroll down to find the cast section. Since the current view
 1304 only shows the movie description and genre tags, I need to
 1305 swipe up to reveal more content, specifically the cast list.

1306 ****Action****
 1307 Swipe up to scroll down the page and reveal more content,
 1308 including the cast section, to find the first actor.
 1309

1310 </think>
 1311
 1312 <Observation>
 1313
 1314 <image_-10-41_1_5.jpg>
 1315

1316 ****Observation****
 1317 1. Open the IMDb app. 2. Search for the movie "Oppenheimer". 3.
 1318 Type "Oppenheimer" into the search bar. 4. Tap on the search
 1319 result for "Oppenheimer" to view its details. 5. Swipe up to
 1320 scroll down the page and reveal more content, including the
 1321 cast section.

1322 </Observation>
 1323

1324 ===Step5===
 1325

1326 <think>
 1327
 1328 <image_-10-41_1_5.jpg>
 1329

1330 ****Thought****
 1331 The current screen shows the movie details for "Oppenheimer". The
 1332 goal is to find the first actor listed in the cast section.
 1333 The "Important Notes" state that the current view does not
 1334 show the cast section, so further scrolling is needed.
 1335 Therefore, the next logical step is to scroll down to reveal
 1336 the cast section.

1337 ****Action Thought****
 1338 The current subgoal is to identify the first actor listed in the
 1339 cast section. Based on the screen information and the previous
 1340 action, I have already swiped up to reveal more content.
 1341 However, the cast section is still not visible in the current
 1342 screenshot. Therefore, I need to swipe up again to scroll
 1343 further down the page until the cast section appears.

1344 ****Action****
 1345 Swipe up again to scroll down the page and reveal more content,
 1346 specifically the cast section, to find the first actor.
 1347

1348 </think>
 1349 <Observation>

```

1350
1351 <image_-10-41_1_6.jpg>
1352
1353 **Observation**
1354 1. Open the IMDb app. 2. Search for the movie "Oppenheimer". 3.
1355 Type "Oppenheimer" into the search bar. 4. Tap on the search
1356 result for "Oppenheimer" to view its details. 5. Swipe up to
1357 scroll down the page and reveal more content, including the
1358 cast section. 6. Identify the first actor listed in the cast
1359 section.
1360 </Observation>
1361
1362 ===Step6===
1363
1364 <think>
1365
1366 <image_-10-41_1_6.jpg>
1367
1368 **Thought**
1369 The current screen shows the movie details page for "Oppenheimer".
1370 The "Cast" section is visible, and the goal is to identify
1371 the first actor listed. The image shows "TOP CAST" and
1372 partially visible images of actors. To identify the first
1373 actor, I need to look at the first visible actor's image and
1374 name. Since the image is only partially visible, I need to
1375 scroll to the right to fully see the first actor's image and
1376 name.
1377 **Action Thought**
1378 The current subgoal is to identify the first actor listed in the
1379 cast section. The screenshot shows the "Cast" section with "
1380 TOP CAST" visible. Below "TOP CAST", there are images of
1381 actors. The first image is partially visible. To identify the
1382 first actor, I need to see the name associated with the first
1383 image. Since the name is not visible, I need to scroll down
1384 further to reveal the full cast list or at least the names of
1385 the top cast members.
1386 **Action**
1387 Swipe up to scroll down the page and reveal the names of the
1388 actors in the "TOP CAST" section.
1389
1390 </think>
1391
1392 <Observation>
1393
1394 <image_-10-41_1_7.jpg>
1395
1396 **Observation**
1397 1. Open the IMDb app. 2. Search for the movie "Oppenheimer". 3.
1398 Type "Oppenheimer" into the search bar. 4. Tap on the search
1399 result for "Oppenheimer" to view its details. 5. Identify the
1400 first actor listed in the cast section.
1401 </Observation>"
1402     }
1403 ],
    "solution": "<answer> Cillian Murphy </answer>"

```

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

}

Listing 2: Example of Execution Trajectory 2

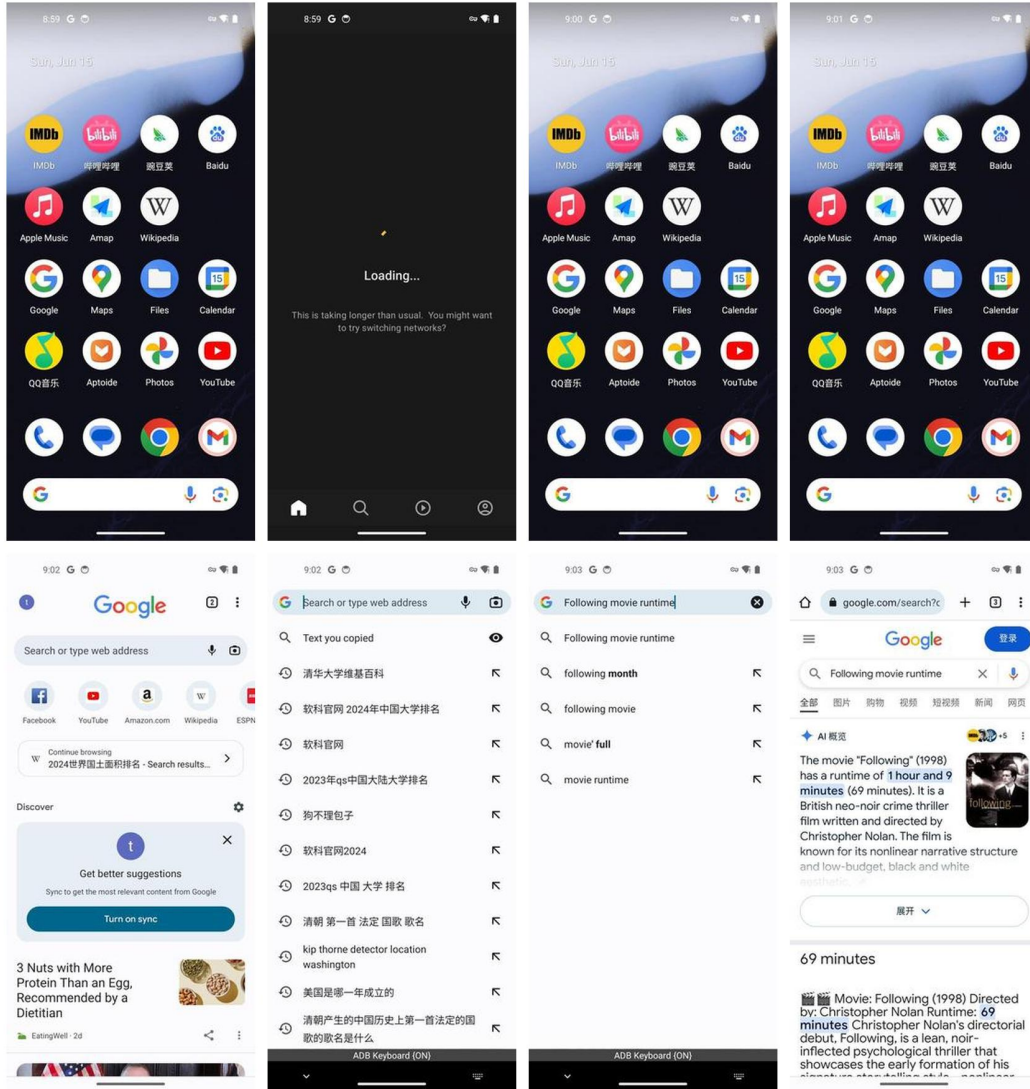


Figure 6: Image from image_59-16.1.1 . jpg to image_59-16.1.8 . jpg

```
{
  "images": [
    "image_16-59-16.1.1 . jpg",
    "image_16-59-16.1.2 . jpg",
    "image_16-59-16.1.3 . jpg",
    "image_16-59-16.1.4 . jpg",
    "image_16-59-16.1.5 . jpg",
    "image_16-59-16.1.6 . jpg",
    "image_16-59-16.1.7 . jpg",
    "image_16-59-16.1.8 . jpg"
  ],
  "messages": [
    {
```

```

1458         "role": "user",
1459         "content": "Please use IMDb to search for the movie '
1460             Following', and tell me what its runtime is?"
1461     },
1462     {
1463         "role": "assistant",
1464         "content": "===Step1===
1465
1466 <think>
1467 <image_16-59-16_1_1.jpg>
1468
1469 **Thought**
1470 The user wants to find the runtime of the movie "Following". I can
1471 use a search engine or a movie database app to find this
1472 information. Since there is an IMDb app on the home screen, it
1473 would be the most direct way to get movie information. If
1474 IMDb doesn't work, I can use Google or Baidu to search for it.
1475
1476 **Action Thought**
1477 The user wants to search for movie information, and the overall
1478 plan is to use the IMDb app. The current subgoal is to open
1479 the IMDb app. Looking at the screen information, I can see an
1480 icon labeled "IMDb" at coordinates [162, 570]. This is the
1481 correct app to open.
1482
1483 **Action**
1484 Open the IMDb app to start searching for the movie information.
1485 </think>
1486 <Observation>
1487 <image_16-59-16_1_2.jpg>
1488
1489 **Observation**
1490 1. Open the IMDb app.
1491
1492 </Observation>
1493
1494 ===Step2===
1495
1496 <think>
1497 <image_16-59-16_1_2.jpg>
1498
1499
1500 **Thought**
1501 The IMDb app is stuck on a loading screen, indicating a potential
1502 network issue. The message "This is taking longer than usual.
1503 You might want to try switching networks?" suggests that the
1504 app is having trouble connecting to the internet. Since the
1505 current plan relies on the IMDb app, and it's not functioning,
1506 I need to revise the plan to address this. The most
1507 straightforward way to search for movie information online is
1508 to use a web browser. Therefore, I will switch to using the
1509 Chrome browser to search for the movie's runtime.
1510
1511 **Action Thought**

```

1512 The current subgoal is to open the Chrome browser. However, the
1513 previous action was to open the IMDb app, which is currently
1514 stuck on a loading screen. Since the IMDb app is not loading,
1515 it's better to close it and open Chrome instead of waiting. To
1516 do this, I need to go back to the home screen first, then
1517 open Chrome.

1518 ****Action****
1519 Return to the home screen to close the stuck IMDb app and prepare
1520 to open the Chrome browser.
1521

1522 `</think>`
1523

1524 `<Observation>`
1525

1526 `<image_16-59-16_1_3.jpg>`
1527

1528 ****Observation****
1529 1. Open the IMDb app. 2. Open the Chrome browser.
1530

1531 `</Observation>`
1532

1533 `===Step3===`
1534

1535 `<think>`
1536

1537 `<image_16-59-16_1_3.jpg>`
1538

1539 ****Thought****
1540 The current plan is to open the Chrome browser. The previous
1541 attempt to use the IMDb app failed because it was stuck
1542 loading. The Chrome browser icon is visible on the screen. I
1543 need to tap on the Chrome icon to open it.

1544 ****Action Thought****
1545 The current subgoal is to open the Chrome browser. Looking at the
1546 screenshot, I can see the Chrome icon at coordinates [663,
1547 1888]. I will use the 'Open_App' action to open Chrome.

1548 ****Action****
1549 Open the Chrome browser to search for the movie runtime.
1550

1551 `</think>`
1552

1553 `<Observation>`
1554

1555 `<image_16-59-16_1_4.jpg>`
1556

1557 ****Observation****
1558 1. Open the IMDb app. 2. Open the Chrome browser.
1559

1560 `</Observation>`
1561

1562 `===Step4===`
1563

1564 `<think>`
1565

1566 `<image_16-59-16_1_4.jpg>`

1566 ****Thought****
1567 The previous subgoal was to open the Chrome browser. The current
1568 screenshot shows the home screen with the Chrome browser icon
1569 visible. The IMDb app was loading slowly, so the plan was to
1570 switch to Chrome. I need to tap on the Chrome icon to open it.
1571

1572 ****Action Thought****
1573 The previous attempt to open Chrome failed. Looking at the
1574 screenshot, I can see the Chrome icon at coordinates [663,
1575 1888]. I will try to tap on this icon to open the Chrome
1576 browser.

1577 ****Action****
1578 Tap the Chrome browser icon to open it.
1579

1580 </think>
1581
1582 <Observation>
1583
1584 <image_16-59-16_1_5.jpg>
1585

1586 ****Observation****
1587 1. Open the IMDb app. 2. Open the Chrome browser. 3. Search for
1588 the movie "Following" and its length.
1589

1589 </Observation>
1590

1591 ===Step5===
1592

1593 <think>
1594
1595 <image_16-59-16_1_5.jpg>
1596

1597 ****Thought****
1598 The current status shows that the Chrome browser is open and
1599 displaying the Google search page. The previous subgoal was to
1600 open the Chrome browser, which has been completed. The next
1601 step in the plan is to search for "Following movie runtime". I
1602 need to type this into the search bar.

1603 ****Action Thought****
1604 The current screen is the Chrome browser, and the goal is to
1605 search for "Following movie runtime". The search bar is
1606 visible with the text "Search or type web address". I need to
1607 tap on this search bar to activate the keyboard and then type
1608 the search query.

1609 ****Action****
1610 Tap the search bar to activate the keyboard and prepare to type
1611 the search query.
1612

1613 </think>
1614
1615 <Observation>
1616
1617 <image_16-59-16_1_6.jpg>
1618

1618 ****Observation****
1619

1620 1. Open the IMDb app. 2. Open the Chrome browser. 3. Search for
 1621 the movie "Following" and its length. 4. The keyboard is
 1622 activated.

1623

1624 </Observation>

1625

1626 ===Step6===

1627

1628 <think>

1629

1630 <image_16-59-16_1_6.jpg>

1631

1632 ****Thought****

1633 The current status shows that the Chrome browser is open and the
 1634 search bar is visible with the text "Search or type web
 1635 address". The keyboard is activated. The next step is to type
 1636 "Following movie runtime" into the search bar.

1637 ****Action Thought****

1638 The current subgoal is to search for "Following movie runtime".
 1639 The Chrome browser is open, and the keyboard is activated. The
 1640 search bar is visible and ready for input. Therefore, the
 1641 next logical step is to type the search query "Following movie
 1642 runtime" into the search bar.

1643 ****Action****

1644 Type "Following movie runtime" into the search bar to search for
 1645 the movie's runtime.

1646

1647 </think>

1648

1649 <Observation>

1650

1651 <image_16-59-16_1_7.jpg>

1652

1653 ****Observation****

1654 1. Open the IMDb app. 2. Open the Chrome browser. 3. Search for
 1655 the movie "Following" and its length. 4. The keyboard is
 1656 activated. 5. Typed "Following movie runtime" into the search
 1657 bar.

1658 </Observation>

1659

1660 ===Step7===

1661

1662 <think>

1663

1664 <image_16-59-16_1_7.jpg>

1665

1666 ****Thought****

1667 The current status shows that "Following movie runtime" has been
 1668 typed into the search bar and the first search suggestion is "
 1669 Following movie runtime". To proceed with the search, I need
 1670 to tap on the search button or the first search suggestion.
 1671 Since the first suggestion is exactly what I typed, tapping on
 1672 it will initiate the search.

1673 ****Action Thought****

```

1674 The current screen shows the Google search page with the search
1675 query "Following movie runtime" already typed in the search
1676 bar. There are also search suggestions below. The first
1677 suggestion is "Following movie runtime", which is exactly what
1678 we want to search for. To initiate the search, I need to
1679 either tap on the search suggestion or press the Enter key.
1680 Tapping the first search suggestion is a more direct way to
1681 proceed.
1682
1683 **Action**
1684 Tap the first search suggestion "Following movie runtime" to
1685 initiate the search.
1686
1687 </think>
1688 <Observation>
1689
1690 <image_16-59-16_1_8.jpg>
1691
1692 **Observation**
1693 1. Open the IMDb app. 2. Open the Chrome browser. 3. Search for
1694 the movie "Following" and its length. 4. The keyboard is
1695 activated. 5. Typed "Following movie runtime" into the search
1696 bar. 6. The search results for "Following movie runtime" are
1697 displayed, and the movie's runtime is visible.
1698
1699 </Observation>"
1700     },
1701     "solution": "<answer> 69 minutes </answer>"
1702 }

```

Listing 3: Example of Execution Trajectory 3

F TRAINING DETAILS

The model was trained using an iterative Reinforcement Fine-Tuning (RFT) methodology with Qwen2.5-VL-7B-Instruct as the base model. The training was conducted on a computing cluster equipped with four NVIDIA A100 GPUs, running Ubuntu 20.04 and utilizing the ms-swift framework. The RFT process spanned three iterations. **This process was initiated with a preliminary version of the data, consisting of 278 human-verified trajectories.** In each iteration, the model generated 16 candidate trajectories for each sample from this progressively augmented dataset. These candidates were subsequently scored and filtered by a Process Reward Model (PRM, qwen-max), with only the high-quality trajectories being retained for the fine-tuning step. The model was then fine-tuned on this curated data for a single epoch using Low-Rank Adaptation (LoRA). Key hyperparameters for the LoRA training included a rank of 8, an alpha of 32, and application to all linear modules. A learning rate of 1×10^{-4} was used with a cosine schedule and a warmup ratio of 0.05. The training was configured with a global batch size of 128 (achieved with a per-device batch size of 1 and 32 gradient accumulation steps) and a maximum sequence length of 16,384 tokens. The entire process was optimized using bfloat16 mixed-precision, DeepSpeed ZeRO stage 2, and FlashAttention. The parameters of the visual encoder remained frozen throughout all training iterations.

A APPENDIX

You may include other additional sections here.