

AGUVIS: UNIFIED PURE VISION AGENTS FOR AUTONOMOUS GUI INTERACTION

Yiheng Xu^{*♦} Zekun Wang^{*♦} Junli Wang^{*△} Dunjie Lu[♦]
 Tianbao Xie[♦] Amrita Saha[♦] Doyen Sahoo[♦] Tao Yu^{†♦} Caiming Xiong^{†♣}

[♦]University of Hong Kong [♣]Salesforce Research

[♦]{yhxu,tyu}@cs.hku.hk [♣]cexiong@salesforce.com

<https://aguvvis-project.github.io>

ABSTRACT

Graphical User Interfaces (GUIs) are critical to human-computer interaction, yet automating GUI tasks remains challenging due to the complexity and variability of visual environments. Existing approaches often rely on textual representations of GUIs, which introduce limitations in generalization, efficiency, and scalability. In this paper, we introduce AGUVIS, a unified pure vision-based framework for autonomous GUI agents that operates across various platforms. Our approach leverages image-based observations, and grounding instructions in natural language to visual elements, and employs a consistent action space to ensure cross-platform generalization. To address the limitations of previous work, we integrate explicit planning and reasoning within the model, enhancing its ability to autonomously navigate and interact with complex digital environments. We construct a large-scale dataset of GUI agent trajectories, incorporating multimodal reasoning and grounding, and employ a two-stage training pipeline that first focuses on general GUI grounding, followed by planning and reasoning. Through comprehensive experiments, we demonstrate that AGUVIS surpasses previous state-of-the-art methods in both offline and real-world online scenarios, achieving, to our knowledge, the first fully autonomous pure vision GUI agent capable of performing tasks independently without collaboration with external closed-source models. We will open-source all datasets, models, and training recipes to facilitate future research.

1 INTRODUCTION

Graphical User Interfaces (GUIs) are a cornerstone of human-computer interaction, providing a structured yet intuitive platform for users to accomplish tasks across various digital environments: website, desktop, and mobile devices (Deng et al., 2023; Zhou et al., 2024; Xie et al., 2024; Rawles et al., 2024b). Automating GUI operations through autonomous agents can revolutionize productivity by enabling seamless task execution on various applications using existing human-centric tools. Moreover, this approach lays the groundwork for advanced AI systems that can interact with and learn from rich digital environments in ways that mirror human behavior.

To effectively perform GUI tasks autonomously, a GUI agent requires three core competencies: understanding, grounding, and planning & reasoning. For GUI understanding, the agent must first comprehend high-resolution and complex interfaces designed for human users, enabling it to grasp the context and perform subsequent reasoning tasks. GUI grounding involves mapping natural language instructions to visual observations of the interface. For planning and reasoning, the agent must synthesize and analyze the current multimodal observations of the environment with previous observations and action histories, enabling it to generate coherent and effective next steps to ultimately achieve the task goal. Although recent advances in large vision-language models (LVLMs) (OpenAI, 2024; Reid et al., 2024; Li et al., 2024a; Wang et al., 2024a) have significantly enhanced the ability of AI systems to interpret complex visual interfaces, there remain critical challenges in grounding and reasoning specifically tailored for GUI tasks. We identify three primary challenges that must be addressed to advance the capabilities of GUI agents:

^{*}Equal contribution [†]Corresponding authors. Work was partially done during YX’s internship at Salesforce.

Enhancing Pure Vision Framework. Previous approaches (Gur et al., 2024; Kim et al., 2023; Deng et al., 2023; Zhou et al., 2024; Xie et al., 2024) predominantly focus on mapping natural language instructions to textual representations of GUIs, such as HTML or accessibility trees. This method presents several limitations. Firstly, GUIs are inherently visual, and leveraging image-based representations aligns more closely with human cognitive processes. Secondly, textual representations can vary widely across different environments, complicating the generalization of the model and limiting the availability of consistent training data. Finally, these textual representations are often verbose and complex, leading to increased inference times compared to more compact image encodings (Figure 2). By unifying observations across platforms as images and grounding instructions to image coordinates, GUI agents can generalize more effectively across diverse environments.

Unification Across GUI Environments. The action spaces and control APIs for GUI interactions vary significantly across diverse environments, particularly when the observations are textual. Even within the same platform, the action space can differ greatly. This heterogeneity limits the amount of training data available for each environment, impeding the development of a model that can generalize effectively across different platforms and scale further. A unified action space that abstracts these environmental differences is crucial for creating robust and adaptable GUI agents. Previous work (Chen et al., 2024b; Zeng et al., 2024) has attempted to unify digital agent data across diverse environments, such as combining GUI, game, and CLI interfaces for joint training. However, these interfaces do not share the same interaction logic. In contrast, GUIs on desktop, web, and mobile platforms naturally share similar human-computer interaction (HCI) logic. This commonality facilitates their unification, enabling consistent visual observations and action spaces that mutually benefit both visual grounding and reasoning.

Integrating Planning and Reasoning with Grounding. Current methodologies (Zheng et al., 2024a) often depend on the reasoning capabilities of closed-source large language models (LLMs) (OpenAI, 2024) to plan the completion of GUI tasks or, alternatively, train agents to make direct action decisions through grounding without an explicit reasoning process. This dichotomy results in either a lack of grounding abilities or a lack of comprehensive reasoning abilities. Recently, some works (Gou et al., 2024; Lu et al., 2024) attempt to use closed-source LLMs with specialized GUI grounding models together and communicate with natural language instruction to utilize both abilities. However, on the one hand, natural language communication between the two models usually results in information loss. On the other hand, most importantly, this approach is not further scalable to solve GUI interaction since grounding has been improved close to the upper bound with data synthesis, and most remaining problems are planning related. However, the GUI planning and reasoning ability of closed-source LLMs cannot be further improved.

To address these challenges, we introduce a unified framework for GUI agents that harmonizes pure vision observation and consistent action spaces across diverse environments. Our approach leverages vision-based grounding to improve generalization and reduce inference costs while employing a standardized action space with a plugin system to facilitate consistent learning and interaction across various platforms. After a unified GUI grounding training stage, we demonstrate that unified augmented datasets can effectively build a model capable of executing complex GUI grounding instructions on various platforms. In addition, we integrate explicit visual planning and reasoning into the same model, enabling autonomous navigation and interaction within complex digital environments. Since existing GUI agent trajectories do not fully support these demands, we have unified the existing planning datasets on different platforms and constructed a large-scale, pure vision, cross-platform, multi-step dataset of agent trajectories, featuring comprehensive multimodal reasoning and grounding. Through extensive experiments across various scenarios, we demonstrate the effectiveness of our approach in advancing the state-of-the-art for pure vision-based autonomous GUI agents. To our knowledge, this is the first model that can autonomously complete tasks in real-world online environments without relying on higher reasoning abilities from closed-source models.

Our contributions are as follows:

- We introduce a unified pure vision framework for building generalizable GUI agents that operate with vision-based observations and a plugin-enabled action system, enhancing cross-platform adaptability.
- We develop a comprehensive data pipeline that unifies existing GUI grounding annotations and integrates explicit planning and reasoning. This enables the construction of large-scale datasets for grounding and multi-step agent trajectory datasets across platforms.

mobile. Over this action space, an agent model can then learn to generate actions in order to control GUI without any action space description.

While mouse and keyboard inputs form the core of GUI interactions, they are not comprehensive. Certain platforms require additional actions. For example: (1) specific actions on mobile platforms such as swiping; (2) shortcuts that efficiently perform a series of actions like opening apps; (3) communication actions such as providing answers or terminating after completion. To address these extended requirements, we introduce a pluggable action system. This system allows us to expand the action space by aligning new actions with the existing `pyautogui` commands where possible. For actions that cannot be directly mapped, the pluggable system provides the flexibility to incorporate them with detailed action descriptions. This enables the model to generalize effectively to environments where new actions are introduced. By combining pure vision observations with a unified action space and a flexible pluggable system, our framework enables the training of a single model that can operate across diverse platforms. This setup not only simplifies the training process but also ensures the model can generalize and adapt to novel environments and tasks.

2.3 THE AGUVIS COLLECTION

GUI agent trajectories are a low-resource data source compared with its challenges. This is because the observation and action space vary across different environments even on the same platform. Fortunately, GUI environments share the same operation logic and similar action space. We can efficiently unify existing data to scale the training set. Therefore, we propose THE AGUVIS COLLECTION, a large-scale GUI agent training dataset collected and augmented with existing GUI agent data. This data collection consists of two splits: grounding split (Table 10) and planning & reasoning split (Table 11), corresponding to the two important GUI abilities.

Template-augmented Grounding Data. Vision-based grounding requires the model to ground the natural language intent to the image observation with coordinates. On one hand, there are several previous works that have built datasets on different platforms, including natural language instructions and corresponding target elements. We collected and unified them into `pyautogui` commands format. On the other hand, we found that there are many datasets proposed for user interfaces on different platforms that contain a large amount of metadata, including the positions of all text/icons/widgets in the current interface. Using this type of data we constructed templates for `pyautogui` actions. We randomly generated grounding data pairs through these templates to train models to ground these elements based on images. This operation greatly expanded the data scale.

VLM-augmented Planning & Reasoning Trajectories. High-quality GUI agent trajectories contain several key components: a high-level goal, a sequence of interleaved observations, natural language reasoning, and grounded actions. Existing approaches typically rely on human annotation to collect these trajectories (Deng et al., 2023; Rawles et al., 2024b; Li et al., 2024c). Most of the agent trajectory data contains high-level goals, observations, and grounded actions. However, the intermediate reasoning process and low-level action instructions are not included. This makes it difficult for existing data to train agents to perform chain-of-thought or inner monologue reasoning to help the model plan the next action, resulting in poor agent performance.

To augment the agent trajectories with detailed reasoning and low-level action instructions, we employ a vision-language model (VLM) to generate the inner monologue for each step in the trajectory. Specifically, for each time step t , given the high-level goal G , the current image observation o_t , and the grounded action a_t , we prompt the VLM to produce the inner monologue components: observation description d_t , thoughts h_t , and low-level action instruction a_t^{instr} . To assist the VLM in generating accurate and contextually relevant monologues, we highlight the target element associated with the grounded action a_t on the image observation o_t . This visual cue helps the model focus on the relevant part of the interface. Additionally, we include the previous low-level action instructions $a_1^{instr}, a_2^{instr}, \dots, a_{t-1}^{instr}$ to provide the VLM with the action history, ensuring continuity and coherence in the generated reasoning.

The prompting strategy is carefully crafted to guide the VLM in generating inner monologues that are predictive and goal-oriented, without relying on hindsight or revealing future actions. By simulating the agent’s thought process in a first-person perspective, we encourage the generation of

actionable instructions that align with the high-level goal and current observation. This approach results in a large-scale dataset of agent trajectories enriched with detailed reasoning and instructions.

2.4 MODEL ARCHITECTURE

Unlike grounding agents that rely on structured UI representations (such as accessibility trees) as their textual input, vision-based grounding requires the model to map intents directly to visual observations. This means the model needs to encode high-resolution images while preserving their original aspect ratios. Recent advances in VLMs have made these capabilities possible. We choose Qwen2-VL (Wang et al., 2024b) as our starting VLM. It uses NaViT as an image encoder with native dynamic resolution support (Dehghani et al., 2023). Unlike its predecessor, Qwen2-VL can now process images of any resolution, dynamically converting them into a variable number of visual tokens. To support this feature, ViT is modified by removing the original absolute position embeddings and introducing 2D-RoPE (Su et al., 2024) to capture the two-dimensional positional information of images. Based on these unique features, Qwen2-VL is highly suitable for GUI agents’ needs. It can encode high-resolution images of any ratio with relatively fewer image token costs. Therefore, we chose Qwen2-VL as our starting VLM to build our GUI agent.

LLaVA-OneVision (Li et al., 2024a) is another suitable VLM as it also supports high-resolution any ratio image encoding, although its image token cost is relatively higher than Qwen2-VL. We also apply our data recipe and training strategy to LLaVA and show that our framework is model-independent and generally works for high-resolution VLMs details are shown in Section 4.2..

2.5 TRAINING PARADIGM

We begin with a Vision-Language Model (VLM) that possesses advanced image understanding capabilities, and the training process is divided into two main stages: Grounding Training and Planning & Reasoning Training. Each stage utilizes a distinct data split from our THE AGUVIS COLLECTION to progressively enhance the VLM’s abilities.

Stage 1: Grounding Training In this stage, we focus on enabling the model to understand and interact with objects within a single GUI screenshot. GUI environments typically feature multiple interactable objects within a single screenshot, generating a large volume of grounding data but leading to shorter, less diverse interaction sequences, which can limit training efficiency.

We train our model with a grounding packing strategy where multiple instruction-action pairs are bundled into a single image, resulting in a single-image-multiple-turn format. This technique allows the model to process several grounding examples from one screenshot, reducing redundant training overhead while retaining a high level of grounding performance. This approach significantly accelerates training by maximizing the use of each image without compromising accuracy. To equip our model with the capability for GUI understanding and grounding, which serves as the foundation for subsequent planning and reasoning, we conducted this training stage. Upon completing Stage 1 training, the model is referred to as AGUVIS-G.

Stage 2: Planning & Reasoning Training Building on the foundation of AGUVIS-G, the second stage introduces more complex decision-making and reasoning processes. This phase is designed to teach the model how to execute multi-step tasks by reasoning through agent trajectories that vary in complexity and environments, encompassing diverse reasoning modes.

Thanks to our detailed inner monologue trajectory data, we implement a reasoning mixture approach, where the model is exposed to various levels of cognitive complexity, from straightforward low-level action instructions to full inner monologues that include observation descriptions, thoughts, and detailed action plans. By dynamically adjusting the complexity of these trajectories, we train the model to be adaptable, fostering step-by-step reasoning and high-level decision-making abilities. This diversity in reasoning ensures that the model can handle a wide range of tasks with nuanced understanding and precision. After this stage, the fully trained model is called AGUVIS, which can be employed in both offline and online GUI tasks across diverse environments.

As shown in Figure 2, AGUVIS significantly reduces the efficiency costs by reducing 93% USD costs and 70% input tokens per step compared to GPT-4o, which indicates considerable potential for applying purely visual agents in practical applications.

4.4 ERROR ANALYSIS

We conduct an error analysis of AGUVIS on 50 samples from the ScreenSpot dataset under the self-plan setting to understand the impact of planning on performance. As shown in Figure 3, our findings reveal that 40% of errors are due to ambiguous instructions that could refer to multiple grounding targets, while the remaining 60% are grounding errors. We observe that in these error cases, the model tends to perform direct grounding action rather than planning explicitly before acting. Notably, when we enforce planning by prompting the agent model to generate low-level instructions before execution, it resolved 20% of the grounding errors. This suggests that while the agent model possesses strong grounding capabilities, there remains significant potential for improvement in effectively leveraging planning and reasoning. These insights highlight opportunities for future work, including improving instruction clarity through the agent model itself, developing adaptive planning mechanisms, and refining training data to include more diverse planning scenarios. Addressing these aspects could further enhance our GUI agent model’s robustness on various tasks and environments.

5 RELATED WORK

5.1 BENCHMARKS AND DATASETS FOR GUI AGENT

Recent advancements in autonomous GUI agents have led to the development of numerous benchmarks and datasets. Web-based benchmarks such as Mind2Web (Deng et al., 2023), WebArena (Zhou et al., 2024; Koh et al., 2024a), WebLNX (Lù et al., 2024), WorkArena (Drouin et al., 2024) and WebCanvas (Pan et al., 2024b) focus on evaluating agents’ performance in web environments. For desktop and mobile platforms, datasets like OSWorld (Xie et al., 2024), WindowsAgentArena (Bonatti et al., 2024), AitW (Rawles et al., 2024b), AitZ (Zhang et al., 2024b), AMEX (Chai et al., 2024), GUI-Odyssey (Lu et al., 2024) and AndroidControl (Li et al., 2024b) have been introduced to assess agents’ capabilities across different operating systems and device types. Cross-platform datasets such as ScreenSpot (Cheng et al., 2024), OmniACT (Kapoor et al., 2024), GUICourse (Chen et al., 2024a), and CRAB (Xu et al., 2024a) aim to provide comprehensive evaluation frameworks spanning multiple devices and interfaces. Evaluations on specialized applications have also emerged, such as WonderBread (Wornow et al., 2024)’s focus on business process management tasks and Spider-2V (Cao et al., 2024)’s on data science and engineering workflows. In this work, we extensively test benchmarks under both online and offline task settings to thoroughly evaluate and demonstrate the model’s planning and grounding capabilities.

5.2 MODELS AND APPROACHES FOR GUI AGENT

In parallel with dataset development, significant progress has been made in creating more capable GUI agents. Models like WebGPT (Nakano et al., 2021), Lemur (Xu et al., 2024b), Agent-Lumos (Yin et al., 2024), CogAgent (Hong et al., 2024), AutoWebGLM (Lai et al., 2024) and xLAM (Zhang et al., 2024a) have demonstrated improved performance in web navigation tasks. Auto-GUI (Zhang & Zhang, 2024), AppAgent (Zhang et al., 2023), and ScreenAgent (Niu et al., 2024) propose novel approaches for direct GUI interaction without relying on application-specific APIs. SearchAgent (Koh et al., 2024b) introduces an inference-time search algorithm to enhance multi-step reasoning and planning in interactive web environments. These advancements collectively contribute to developing more sophisticated and capable GUI agents, pushing the boundaries of what’s possible in automated task completion across various digital platforms.

6 CONCLUSION

In this paper, we introduced AGUVIS, a unified pure vision-based framework for building autonomous GUI agents that operate across diverse platforms. By only leveraging vision-based observations and a consistent action space, AGUVIS addresses the key challenges of GUI grounding,

Jianguo Zhang, Tian Lan, Ming Zhu, Zuxin Liu, Thai Hoang, Shirley Kokane, Weiran Yao, Juntao Tan, Akshara Prabhakar, Haolin Chen, Zhiwei Liu, Yihao Feng, Tulika Awalgaonkar, Rithesh Murthy, Eric Hu, Zeyuan Chen, Ran Xu, Juan Carlos Niebles, Shelby Heinecke, Huan Wang, Silvio Savarese, and Caiming Xiong. xlam: A family of large action models to empower ai agent systems. 2024a. URL <https://api.semanticscholar.org/CorpusID:272424184>.

Jiwen Zhang, Jihao Wu, Yihua Teng, Minghui Liao, Nuo Xu, Xiao Xiao, Zhongyu Wei, and Duyu Tang. Android in the zoo: Chain-of-action-thought for gui agents. *ArXiv preprint*, 2024b. URL <https://arxiv.org/abs/2403.02713>.

Zhuosheng Zhang and Aston Zhang. You only look at screens: Multimodal chain-of-action agents. In *Findings of the Association for Computational Linguistics*, 2024.

Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. Gpt-4v(ision) is a generalist web agent, if grounded. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*, 2024a.

Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. Synapse: Trajectory-as-exemplar prompting with memory for computer control. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024b.

Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web environment for building autonomous agents. In *International Conference on Learning Representations*, 2024.

Table of Contents in Appendix

A AGUVIS Unified Design	17
A.1 Details of Action Space in AGUVIS	17
A.2 Pluggable Functions: Mobile Environments as An Example	17
B Data Curation of THE AGUVIS COLLECTION	18
B.1 Detailed Source Dataset Statistics	18
B.2 Prompt for Augmenting Planning & Reasoning Trajectories	18
B.3 Human Study on Augmented Data	19
B.3.1 Qualitative Human Study	19
B.3.2 Failure Cases Under Noisy Training Data	19
C AGUVIS Training	21
C.1 Training Example Schema	21
C.2 Training Details	22
D Evaluation Benchmarks	22
D.1 GUI Grounding Evaluation	22
D.2 Offline GUI Agent Evaluation	22
D.3 Online GUI Agent Evaluation	24
D.3.1 Prompts for using GPT-4o as Planning Model	25
E Analysis	28
E.1 Training Ablation	28
E.1.1 Training Strategy Ablation	28
E.1.2 Data Strategy Ablation	29
E.2 Planning Analysis	30
E.2.1 Prompts for self-planning and enforced planning mode.	30
E.2.2 Planning Bonus Examples	31
E.3 AGUVIS Trajectories Examples on Online Evaluation	32
E.3.1 Mind2Web-Live Case: AGUVIS-72B as Planner and Grounder	32
E.3.2 Mind2Web-Live Case: GPT-4o as Planner and AGUVIS-7B as Grounder	33
E.3.3 AndroidWorld Case: AGUVIS-72B as Planner and Grounder	34
E.3.4 AndroidWorld Case: GPT-4o as Planner and AGUVIS-7B as Grounder	35
E.4 Case of AGUVIS Generalization in Real-World Scenarios	36

A AGUVIS UNIFIED DESIGN

A.1 DETAILS OF ACTION SPACE IN AGUVIS

In this section, we introduce our unified action space of our pure vision agent framework AGUVIS. As shown in Table 9, we use default standard pyautogui actions with pluggable actions as the action space of AGUVIS, which ensures the agent model’s universality across environments as well as its flexibility in the specific environment.

Table 9: Default standard pyautogui actions \mathcal{A} with pluggable actions.

Category	Action Space
Basic Actions	pyautogui.moveTo(x, y)
	pyautogui.click(x, y)
	pyautogui.write('text')
	pyautogui.press('enter')
	pyautogui.hotkey('ctrl', 'c')
	pyautogui.scroll(200)
Pluggable Actions	pyautogui.dragTo(x, y)
	browser.select_option(x, y, value)
	mobile.swipe(from, to)
	mobile.home()
	mobile.back()
	mobile.open_app(name)
...	terminate(status)
	answer(text)
...	...

A.2 PLUGGABLE FUNCTIONS: MOBILE ENVIRONMENTS AS AN EXAMPLE

In the mobile environment, we provide the following pluggable functions for Aguvis, along with their corresponding descriptions as shown in Figure A.2.

Pluggable Functions for AGUVIS

You are a GUI agent. You are given a task and a screenshot of the screen. You need to perform a series of pyautogui actions to complete the task.

You have access to the following functions:

```

- {"name": "mobile.home", "description": "Press the home button"}
- {"name": "mobile.back", "description": "Press the back button"}
- {
    "name": "mobile.long_press",
    "description": "Long press on the screen",
    "parameters": {
        "type": "object",
        "properties": {"x": {"type": "number", "description": "The x coordinate of the long press"}, "y": {"type": "number", "description": "The y coordinate of the long press"}},
        "required": ["x", "y"]
    }
}
- {
    "name": "mobile.open_app",
    "description": "Open an app on the device",
    "parameters": {
        "type": "object",
    }
}

```

```

    "properties": {"app_name": {"type": "string",
    "description": "The name of the app to open"}},
    "required": ["app_name"]
}
}
-
{
    "name": "terminate",
    "description": "Terminate the current task and report its completion status",
    "parameters": {
        "type": "object",
        "properties": {"status": {"type": "string", "enum": ["success"]}, "description": "The status of the task"}},
        "required": ["status"]
}
}
-
{
    "name": "answer",
    "description": "Answer a question", "parameters": {
        "type": "object",
        "properties": {"answer": {"type": "string", "description": "The answer to the question"}},
        "required": ["answer"]
}
}
}

```

B DATA CURATION OF THE AGUVIS COLLECTION

B.1 DETAILED SOURCE DATASET STATISTICS

We present the detailed statistical information of all training datasets utilized in both the grounding and planning & reasoning stages. The statistics are shown in Table 10 and Table 11, respectively.

Table 10: The grounding split of THE AGUVIS COLLECTION. Each example in this split consists of a single-step trajectory.

Data source	Platform	Instruction	#Trajectory
SeeClick (Cheng et al., 2024)	Website	Augmented	271K
GUIEnv (Chen et al., 2024a)	Website	Augmented	328K
GUIAct (Chen et al., 2024a)	Website	Original	67K
WebUI (Wu et al., 2023)	Website	Augmented	57K
Widget Captioning (Li et al., 2020b)	Mobile	Original	101K
RicoSCA (Li et al., 2020a)	Mobile	Original	173K
UI RefExp (Bai et al., 2021)	Mobile	Original	16K
RICO Icon (Deka et al., 2017)	Mobile	Augmented	16K
OmniACT (Kapoor et al., 2024)	Desktop & Website	Original	7K
Total			1.036M

B.2 PROMPT FOR AUGMENTING PLANNING & REASONING TRAJECTORIES

Prompt for GPT-4o generating planning & reasoning data

```

Goal: {goal}
Previous Actions: {previous_actions}

```

Given the current screenshot and the next ground truth action labeled as `'{current_action_instruction}'`, the action commands is:

between these extra actions and the overall goal, ultimately resulting in incorrect reasoning and planning.

While these redundant actions do not compromise the trajectory's overall completeness or correctness, they do introduce challenges for the VLM in generating accurate planning.

Goal: Show me the page with information about Wi-Fi subscriptions.

Previous actions:
Step 1: Click on the 'TRAVEL INFO' tab to explore available options for Wi-Fi subscription information.
Step 2: Click on the 'Pre-paid Wi-Fi' link to access Wi-Fi subscription information.

Action: pyautogui.click(x=0.8018, y=0.6183)

Generated Planning

Observation: I am on a page related to United Wi-Fi options, which includes a link for Wi-Fi subscriptions. This is in line with my goal to gather information about Wi-Fi subscriptions for flights.

Thought: Identifying the 'Wi-Fi subscriptions' link on the United Wi-Fi page aligns perfectly with the goal of finding subscription information. I am considering available options and see a dedicated link for subscriptions.

Low-level Instruction: Click on the 'Wi-Fi subscriptions' link to access information about Wi-Fi subscriptions.

Goal: Find the documents required to sell a car in New Jersey for which payments are being paid and it has single person on the title who can come to the store

Previous actions:
Step 1: Click on the 'Sell/Trade' link in the navigation bar to proceed.
Step 2: Select 'New Jersey' from the dropdown menu labeled 'In what state are you selling the car?'.

Action: pyautogui.click(x=0.5938, y=0.8633)

Generated Planning

Observation: The page asks if payments are still being made on the car and if there is more than one person on the title, with the option to select either 'Yes' or 'No' for both questions. Currently, the focus is on confirming that payments are being made, as illustrated by the page.

Thought: To determine the document requirements for selling a car in New Jersey, I need to select options that accurately reflect the current status of the car's payments and ownership. Since payments are being made, I will choose 'Yes' for that question to proceed and assess any further requirements or conditions this might bring into play.

Low-level Instruction: Click on the 'Yes' radio button next to 'Are you making payments on the car?'

Figure 4: Examples of augmented planning and reasoning data generated by GPT-4o. The position of the mouse in the image represents the ground truth click position in the training data.

C AGUVIS TRAINING

C.1 TRAINING EXAMPLE SCHEMA

Training Data Schema of Stage 1 Grounding

Prompt

```
<|im_start|>system  
You are a GUI agent. You are given a task and a screenshot of the  
screen. You need to perform a series of pyautogui actions to  
complete the task.<|im_end|>  
<|im_start|>user  
<|vision_start|><|image_pad|><|vision_end|>  
Please generate the next move according to the ui screenshot,  
instruction and previous actions.  
Instruction: {overall_goal}  
Previous actions: {previous_actions}  
<|im_end|>
```

Generation

```
<|im_start|>assistant<|recipient|>os  
Action: {pyautogui function}  
<|diff_marker|>
```

Training Data Schema of Stage 2 Planning

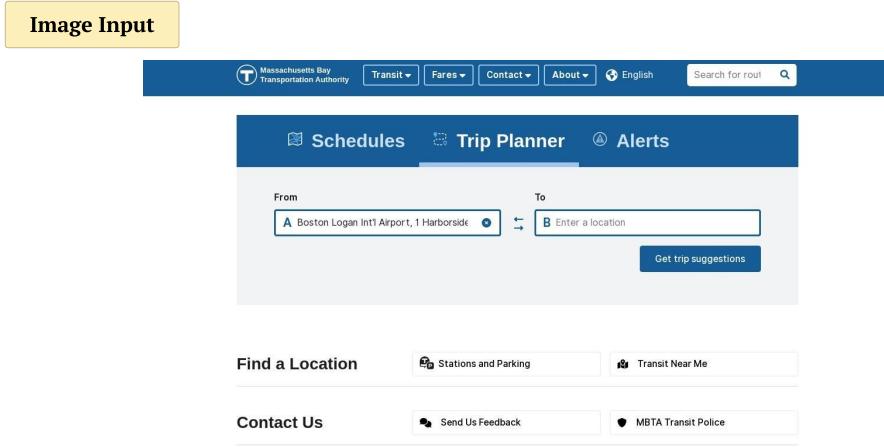
Prompt

```
<|im_start|>system  
You are a GUI agent. You are given a task and a screenshot of the  
screen. You need to perform a series of pyautogui actions to  
complete the task.<|im_end|>  
<|im_start|>user  
<|vision_start|><|image_pad|><|vision_end|>  
Please generate the next move according to the ui screenshot,  
instruction and previous actions.  
Instruction: {overall_goal}  
Previous actions: {previous_actions}  
<|im_end|>
```

Generation

```
<|im_start|>assistant<|recipient|>all  
Observation: {Observation}  
Thought: {Planning}  
Low-level Instruction: {Low-level Instruction}  
<|im_end|>  
<|im_start|>assistant<|recipient|>os  
Action: {pyautogui function}  
<|diff_marker|>
```

AGUVIS introduces a novel explicit planning and reasoning training framework that differs from existing approaches. We illustrate these differences with visual examples in Figure 5. While existing training datasets utilize trajectory data to fine-tune agents, these approaches often involve agents directly outputting action commands (e.g., via pyautogui), bypassing the generation of observations, thoughts, and low-level instructions in natural language that correspond to actions. To elicit the reasoning and planning capabilities of vision-language models and provide the model with richer context for action generation, we scale up training datasets that explicitly require the model to output



Existing GUI Agent Data

Prompt

Please generate the next move according to the UI screenshot, instruction and previous actions.

Instruction: Plan a trip from Boston Logan Airport to North Station.

Previous actions:

```
Step 1: pyautogui.click(x=0.4754, y=0.2062)
Step 2: pyautogui.click(x=0.3295, y=0.4)
pyautogui.write(text='Boston Logan Airport')
Step 3: pyautogui.click(x=0.3262, y=0.4764)
```

Generation

Action:

```
pyautogui.click(x=0.6756, y=0.4)
pyautogui.write(text='North Station')
```

Aguvis Collection Data

Prompt

Please generate the next move according to the UI screenshot, instruction and previous actions.

Instruction: Plan a trip from Boston Logan Airport to North Station.

Previous actions:

```
Step 1: Click on the 'Trip Planner' tab to begin planning the trip.
Step 2: Click on the 'From' input field and type 'Boston Logan Airport'.
Step 3: Click on 'Boston Logan Int'l Airport, 1 Harborside Dr, East Boston, MA 02128, United States' to set it as my starting location.
```

Generation

Observation: The trip planner interface is open with the 'From' field set to Boston Logan Airport. The 'To' field is empty, awaiting input for the destination.

Thought: I have set my starting point as Boston Logan Airport. To proceed, I need to set the destination to North Station, allowing the trip planner to suggest routes.

Low-level Instruction: Click on the 'To' input field and type 'North Station' as the destination.

Action:

```
pyautogui.click(x=0.6756, y=0.4)
pyautogui.write(text='North Station')
```

Figure 5: Compared to the schema of existing gui agent data (left), the schema of AGUVIS planning & reasoning data (right) includes explicit reasoning process with informative natural language previous action context.

AndroidControl. Following the setting in Li et al. (2024d), we randomly sample 500 step-actions from AndroidControl full test set to create a subset, and we report the step accuracy on out-of-domain (OOD) data within both high-level and low-level tasks. The high-level task setting necessitates that the model plans and executes actions, whereas the low-level task setting requires the model to simply adhere to human-labeled instructions for executing the next-step action.

D.3 ONLINE GUI AGENT EVALUATION

Mind2Web-Live. We adopt Mind2Web-Live (Pan et al., 2024b) to evaluate GUI agents’ online planning, a derived dynamic data set from Mind2Web, comprising 104 real-time interactive web tasks. It evaluates whether each required step within a task has been successfully completed and uses the task success rate (Task SR) as the reported metric. The original Mind2Web-Live is built with WebCavas (Pan et al., 2024a), which is a text-based agent framework. To better accommodate the unified observation and action space of pure vision models, we utilize BrowserGym (Drouin et al., 2024) as the evaluation environment for online web tasks which provide support for pure vision-based agent models. BrowerGym is a browser testing environment built on the Playwright (Microsoft, 2024) engine. We incorporate all Mind2Web-Live tasks and evaluation into BrowerGym, involving registering all Mind2Web-Live tasks, setting up the entry points for these tasks, and porting the Mind2Web-Live evaluation functions to BrowerGym.

As Mind2Web-Live is a text-based benchmark, we have to adapt its evaluation function to suit our pure vision-based model. To achieve this, we introduce the two modifications following:

- For the Mind2Web-Live benchmark’s click verification, we adapt our coordinate-based approach by comparing the ground truth CSS selector’s bounding box (when available) with our click coordinates, as we cannot directly identify HTML elements.
- Similarly, for input validation, we retrieve and compare the value of the ground truth input element (if present) with the expected value, circumventing the need for precise HTML element identification based on CSS selectors.

The Mind2Web-Live environment relies on real-world websites, many of which implement detection systems for automated browser testing and reCAPTCHA challenges. These factors created difficulties during evluation on the Mind2Web-Live dataset, resulting in a lower task success rate (Task SR). Specifically, we observed the following websites to have significant issues with automation detection:

- **kohls.** Model using the search functionality on the Kohls website through Playwright directly results in a 502 Bad Gateway error.
- **target.** We are unable to open target’s job website using Playwright due to network connection error.
- **united.** We are unable to open united website using Playwright due to network connection error.

In addition to the websites that were consistently prone to failure, several other sites intermittently blocked our Playwright access during testing. In total, we encountered 18 network errors and 6 reCAPTCHA tasks that the model was unable to complete, preventing our model from scoring on these 24 tasks.

AndroidWorld. AndroidWorld (Rawles et al., 2024b) is a benchmark operating on an Android virtual environment, capable of dynamically instantiating with randomly generated parameters to generate unique tasks for automatic evaluation. It spans 20 real-world applications, encompassing 116 diverse tasks. To assess the pure vision agent models, we follow the instructions in Rawles et al. (2024b), installing a Pixel 6 phone simulator on our computers to serve as the experimental environment. The benchmark incorporates a fully automated task-level evaluation system that automatically assesses whether a state has successfully completed a designated task. The AndroidWorld environment supports optional inputs such as Set-of-Mark (SoM) and textual AXTree information, which most multimodal models currently rely on to complete tasks. However, we solely use raw screenshots as the observation input and restrict the model to coordinate-level actions and basic mobile functions.

Table 12 – Continued from the previous page

Instructions

Review the current state of the page and all other information to find the best possible next action to accomplish your goal. Your answer will be interpreted and executed by a program, make sure to follow the formatting instructions.

Description: Double click an element. The 'element' parameter represents the semantic information of the element you want to double click.

hover(element: str)

Description: Hover over an element. The 'element' parameter represents the semantic information of the element you want to hover over.

keyboard_press(key: str)

Description: Press a combination of keys. Accepts the logical key names that are emitted in the keyboardEvent.key property of the keyboard events: Backquote, Minus, Equal, Backslash, Backspace, Tab, Delete, Escape, ArrowDown, End, Enter, Home, Insert, PageDown, PageUp, ArrowRight, ArrowUp, F1 - F12, Digit0 - Digit9, KeyA - KeyZ, etc. You can alternatively specify a single character you'd like to produce such as "a" or "#". Following modification shortcuts are also supported: Shift, Control, Alt, Meta.

Only a single action can be provided at once. Example:

fill('comment text area', 'This is an example')

Note: you are on mac so you should use Meta instead of Control for Control+C etc.

Table 13: Prompts used for the planning model in **AndroidWorld**, modified from the prompt in (Rawles et al., 2024a)

Instruction

You are an agent who can operate an Android phone on behalf of a user. Based on user's goal/request, you may

- Answer back if the request/goal is a question (or a chat message), like user asks "What is my schedule for today?".
- Complete some tasks described in the requests/goals by performing actions (step by step) on the phone.

When given a user request, you will try to complete it step by step. At each step, you will be given the current screenshot and a history of what you have done (in text). Based on these pieces of information and the goal, you must choose to perform one of the action in the following list (action description followed by the JSON format) by outputting the action in the correct JSON format.

- If you think the task has been completed, finish the task by using the status action with complete as goal_status: '{ "action_type": "status", "goal_status": "complete"}'
- If you think the task is not feasible (including cases like you don't have enough information or can not perform some necessary actions), finish by using the 'status' action with infeasible as goal_status: '{ "action_type": "status", "goal_status": "infeasible"}'
- Answer user's question: '{ "action_type": "answer", "text": "answer_text"}'
- Click/tap on an element on the screen. Please describe the element you want to click using natural language. '{ "action_type": "click", "target": target_element_description}'.
- Long press on an element on the screen, similar with the click action above, use the semantic description to indicate the element you want to long press: '{ "action_type": "long_press", "target": target_element_description}'.
- Type text into a text field (this action contains clicking the text field, typing in the text and pressing the enter, so no need to click on the target field to start), use the semantic description to indicate the target text field: '{ "action_type": "input_text", "text": text.input, "target": target_element_description}'

Continued on the next page

Table 13 – Continued from the previous page

-
- Press the Enter key: ‘{“action_type”: “keyboard_enter”}’
 - Navigate to the home screen: ‘{“action_type”: “navigate_home”}’
 - Navigate back: ‘{“action_type”: “navigate_back”}’
 - Scroll the screen or a scrollable UI element in one of the four directions, use the same semantic description as above if you want to scroll a specific UI element, leave it empty when scroll the whole screen: ‘{“action_type”: “scroll”, “direction”: up, down, left, right, “element”: optional_target_element_description}’
 - Open an app (nothing will happen if the app is not installed): ‘{“action_type”: “open_app”, “app_name”: name}’
 - Wait for the screen to update: ‘{“action_type”: “wait”}’
-

Guidelines

Here are some useful guidelines you need to follow:

General:

- Usually there will be multiple ways to complete a task, pick the easiest one. Also when something does not work as expected (due to various reasons), sometimes a simple retry can solve the problem, but if it doesn't (you can see that from the history), SWITCH to other solutions.
- Sometimes you may need to navigate the phone to gather information needed to complete the task, for example if user asks "what is my schedule tomorrow", then you may want to open the calendar app (using the ‘open_app’ action), look up information there, answer user's question (using the ‘answer’ action) and finish (using the ‘status’ action with complete as goal_status).
- For requests that are questions (or chat messages), remember to use the ‘answer’ action to reply to user explicitly before finish! Merely displaying the answer on the screen is NOT sufficient (unless the goal is something like "show me ...").
- If the desired state is already achieved (e.g., enabling Wi-Fi when it's already on), you can just complete the task.

Action Related:

- Use the ‘open_app’ action whenever you want to open an app (nothing will happen if the app is not installed), do not use the app drawer to open an app unless all other ways have failed.
- Use the ‘input_text’ action whenever you want to type something (including password) instead of clicking characters on the keyboard one by one. Sometimes there is some default text in the text field you want to type in, remember to delete them before typing.
- For ‘click’, ‘long_press’ and ‘input_text’, the target_element_description parameter you choose must based on a VISIBLE element in the screenshot.
- Consider exploring the screen by using the ‘scroll’ action with different directions to reveal additional content.
- The direction parameter for the ‘scroll’ action can be confusing sometimes as it's opposite to swipe, for example, to view content at the bottom, the ‘scroll’ direction should be set to “down”. It has been observed that you have difficulties in choosing the correct direction, so if one does not work, try the opposite as well.

Text Related Operations:

- Normally to select certain text on the screen: (i) Enter text selection mode by long pressing the area where the text is, then some of the words near the long press point will be selected (highlighted with two pointers indicating the range) and usually a text selection bar will also appear with options like ‘copy’, ‘paste’, ‘select all’, etc. (ii) Select the exact text you need. Usually the text selected from the previous step is NOT the one you want, you need to adjust the range by dragging the two pointers. If you want to select all text in the text field, simply click the ‘select all’ button in the bar.
- At this point, you don't have the ability to drag something around the screen, so in general you can not select arbitrary text.

Continued on the next page

Table 13 – Continued from the previous page

-
- To delete some text: the most traditional way is to place the cursor at the right place and use the backspace button in the keyboard to delete the characters one by one (can long press the backspace to accelerate if there are many to delete). Another approach is to first select the text you want to delete, then click the backspace button in the keyboard.
 - To copy some text: first select the exact text you want to copy, which usually also brings up the text selection bar, then click the ‘copy’ button in bar.
 - To paste text into a text box, first long press the text box, then usually the text selection bar will appear with a ‘paste’ button in it.
 - When typing into a text field, sometimes an auto-complete dropdown list will appear. This usually indicating this is a enum field and you should try to select the best match by clicking the corresponding one in the list.
-

E ANALYSIS

E.1 TRAINING ABLATION

E.1.1 TRAINING STRATEGY ABLATION

To further demonstrate the contribution of Stage 1, Stage 2, and their combination to model training, we conducted an ablation study. Specifically, we designed five experimental settings on AGUVIS_{QWEN2-VL} and AGUVIS_{LLAVA-OV}:

- **Stage 1 → Stage 2** corresponds to the staged configuration AGUVIS used in our paper, where Stage 1 is followed by Stage 2 sequentially.
- **Stage 1 + Stage 2** represents a joint training setup, where two stages are combined into a training process.
- **w/o Stage x** indicates the absence of the respective stage in the setting.

Note that for the setting of removing Stage 2 (w/o Stage 2 or w/o Stage 1 & 2), the models are fine-tuned on the corresponding task-specific dataset for planning tasks.

From the first two rows in Table 14, it can be observed that the differences between models trained with Staged Training and Joint Training setups are relatively minor. However, a clear trend emerges: models trained using the Joint Training setup perform better on GUI grounding tasks but exhibit inferior performance on datasets requires planning ability such as MM-Mind2Web and AndroidControl High-level. This trend implies grounding data in Stage 1 is more abundant, dominating the optimization process and biasing the model toward grounding tasks. In contrast, the data in Stage 2, which combines planning and grounding, is of higher quality and better aligned with the agent’s deployment scenarios. This rationale underpins our decision to position Stage 2 later in the training sequence.

Moreover, it is observed that compared to AGUVIS_{QWEN2-VL} trained through both Stage 1 and Stage 2, the model trained with only Stage 2 data maintains similar performance on MM-Mind2Web and AndroidControl but exhibits a notable decline in GUI grounding performance on ScreenSpot. This suggests that the stability on Mind2Web and AndroidControl can be attributed to Qwen2VL’s pre-training on natural image grounding. However, the diverse image and domain requirements of the ScreenSpot GUI grounding test set highlight the necessity of extensive and varied grounding training from Stage 1. This training is essential for improving the grounding performance required for a cross-platform GUI agent model.

To verify this analysis, we conduct the same ablation study on the LLava model, as shown in Table 15. From the results, we can see that the original LLava did not undergo extensive natural image grounding training during the training process, making it insufficient for LLava to excel when only Stage 1 or Stage 2 is conducted. When both Stage 1 and Stage 2 are performed, LLava can be significantly improved, even surpassing previous SOTA results. This validates the above analysis and further demonstrates that our method is model-agnostic and universally applicable to popular VLMs like Qwen2-VL and LLava.

for subsequent steps. Additionally, incorporating low-level instructions from the training data improves the accuracy of the model’s action execution, as demonstrated in both the Screenspot and AndroidControl low-level tasks.

E.2 PLANNING ANALYSIS

E.2.1 PROMPTS FOR SELF-PLANNING AND ENFORCED PLANNING MODE.

In Appendix C.1, we present the training data schema for Stage 1 and Stage 2. We use the special token <|recipient|> along with os or all to control whether the message content is an inner monologue or a pyautogui action command. Thanks to this design, we can use <|recipient|> during the inference phase to control the content generated by the model.

In the Enforced Plan Setting, we employ the <|recipient|>all\nThought prompt to compel the model to generate a planning phase following this. While in the self-plan setting, we do not add any word after <|recipient|>, so the model can choose to generate os to directly produce a pyautogui command, or generate all to first create natural language reasoning and then generate a pyautogui command.

Prompt Template For Self-plan

```
<|im_start|>system
You are a GUI agent. You are given a task and a screenshot of the
screen. You need to perform a series of pyautogui actions to
complete the task.<|im_end|>
<|im_start|>user
<|vision_start|><|image_pad|><|vision_end|>Please generate the next
move according to the ui screenshot, instruction and previous
actions.

Instruction: {goal}

Previous actions: {previous_actions}
<|im_end|>
<|im_start|>assistant<|recipient|>
```

Prompt Template For Enforced Plan

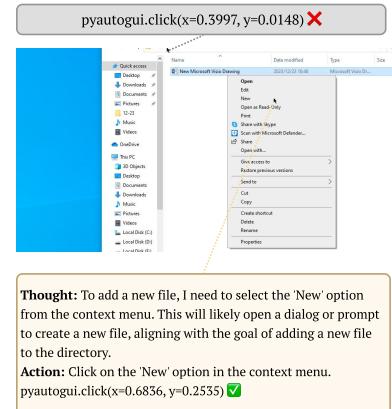
```
<|im_start|>system
You are a GUI agent. You are given a task and a screenshot of the
screen. You need to perform a series of pyautogui actions to
complete the task.<|im_end|>
<|im_start|>user
<|vision_start|><|image_pad|><|vision_end|>Please generate the next
move according to the ui screenshot, instruction and previous
actions.

Instruction: {overall_goal}

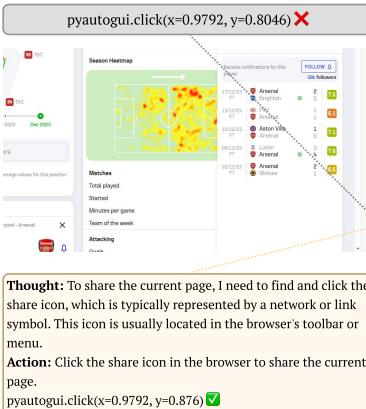
Previous actions: {previous_actions}
<|im_end|>
<|im_start|>assistant<|recipient|>all
Thought:
```

E.2.2 PLANNING BONUS EXAMPLES

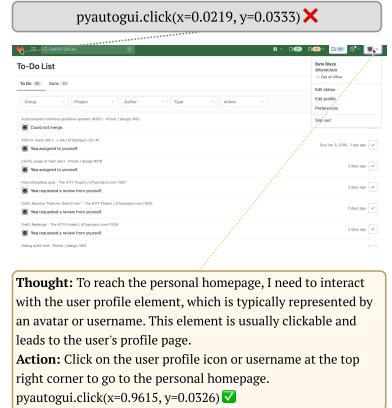
Task Instruction: add a new file



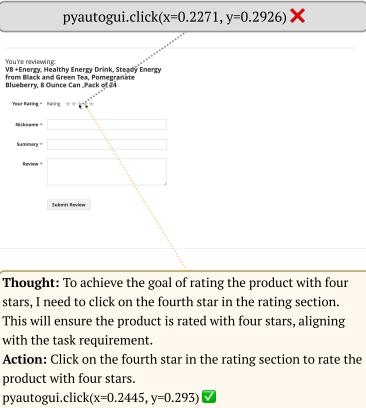
Task Instruction: send current page



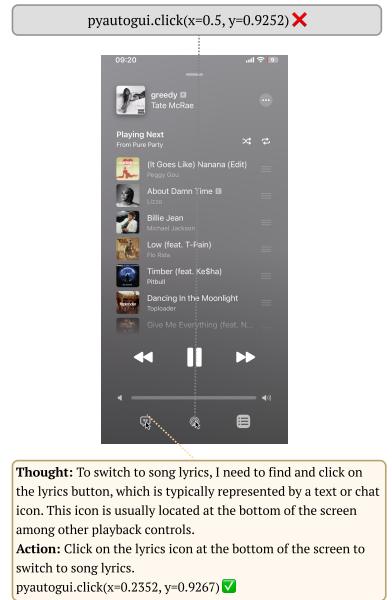
Task Instruction: go to personal homepage



Task Instruction: rate this product with four stars



Task Instruction: switch to song lyric



Task Instruction: change dark mode options

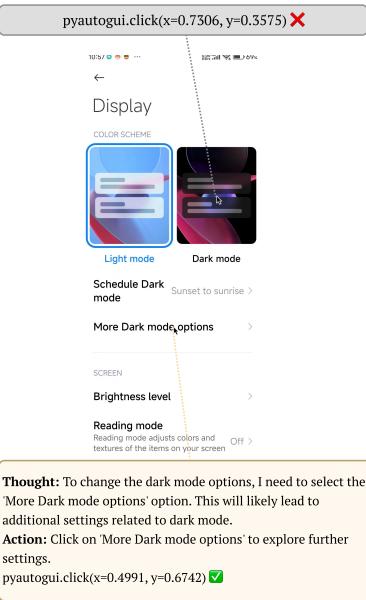


Figure 6: Self-plan examples on different environments.

E.3 AGUVIS TRAJECTORIES EXAMPLES ON ONLINE EVALUATION

E.3.1 MIND2WEB-LIVE CASE: AGUVIS-72B AS PLANNER AND GROUNDER

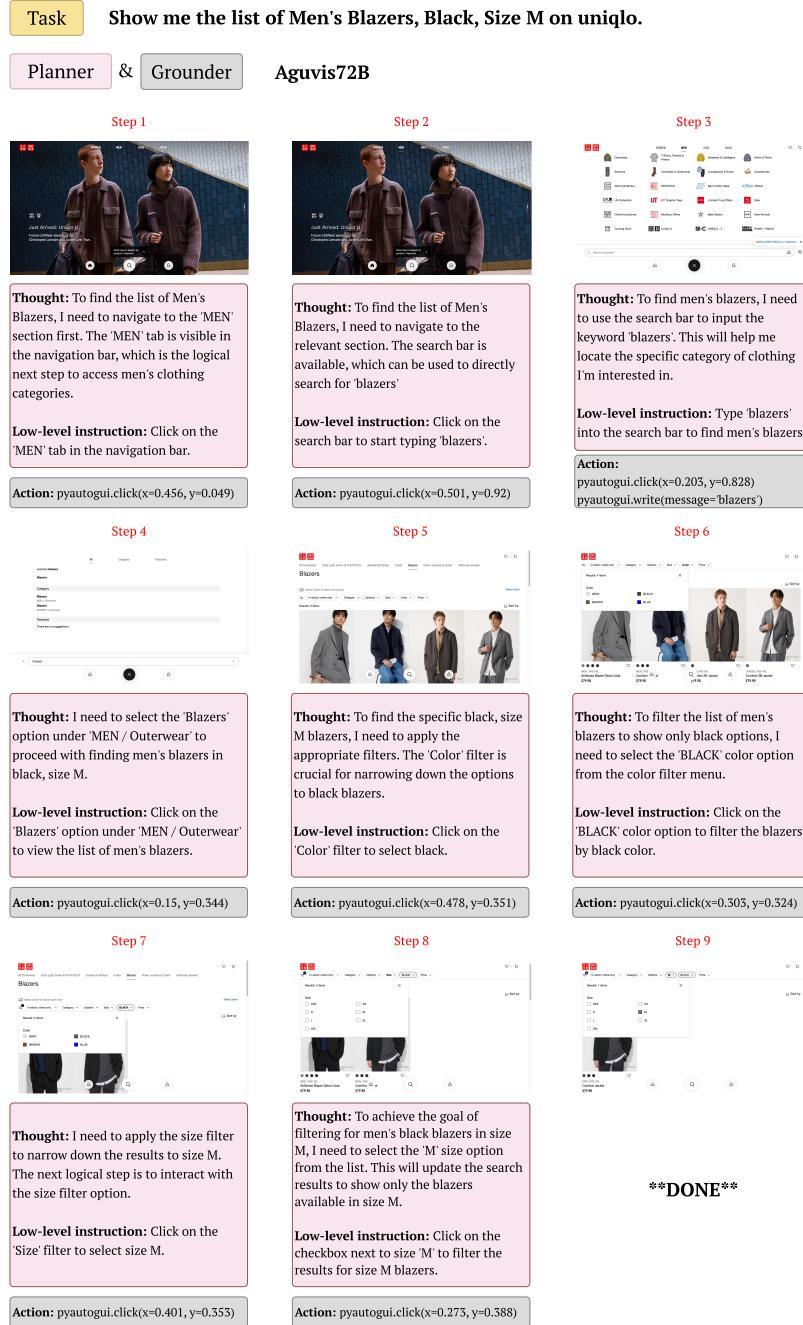


Figure 7: Example of AGUVIS-72B as planner and grounder executing Mind2Web-Live task. Due to space limitations, we present here the trajectory generated guided by Thought.

E.3.2 MIND2WEB-LIVE CASE: GPT-4O AS PLANNER AND AGUVIS-7B AS GROUNDER

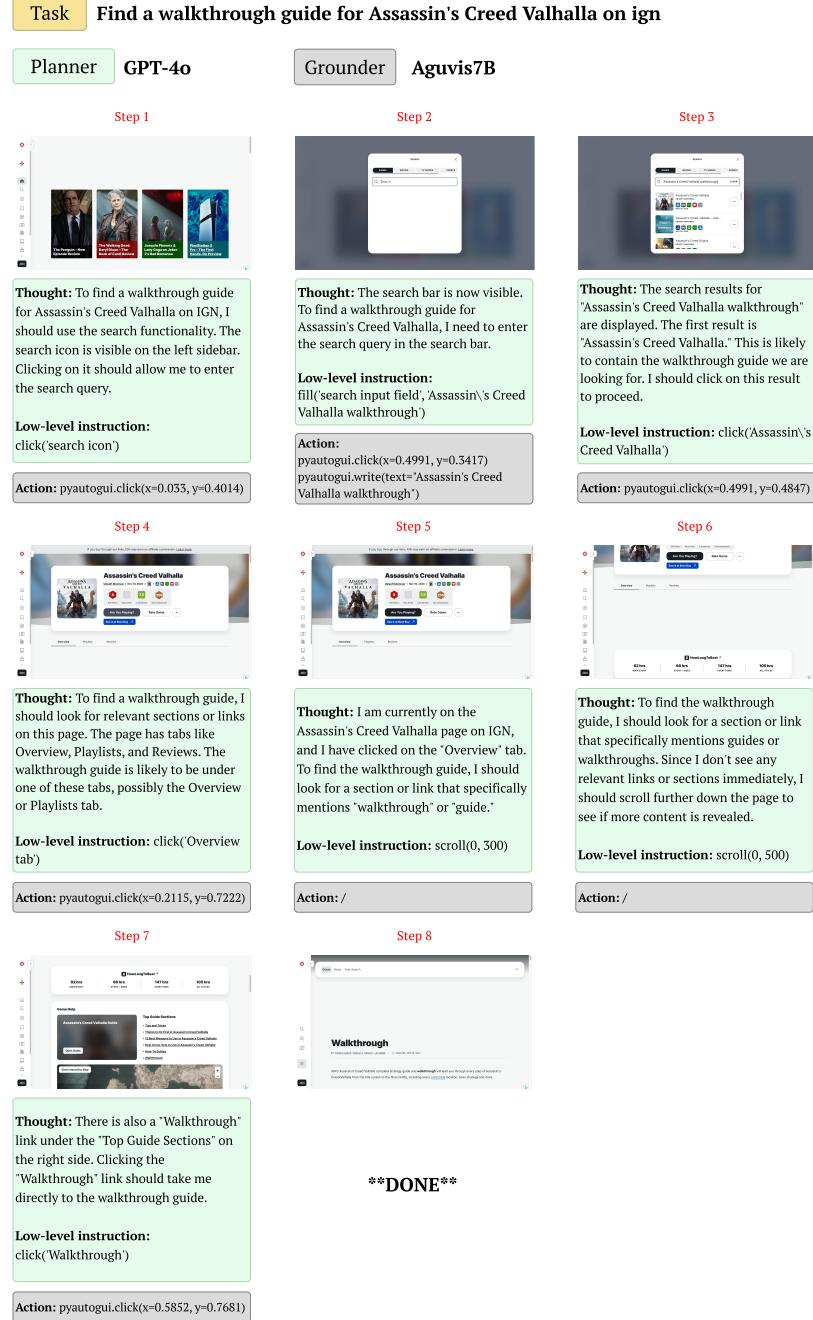


Figure 8: Example of GPT-4o as planner and AGUVIS-7B as grounder executing Mind2Web-Live task.

E.3.3 ANDROIDWORLD CASE: AGUVIS-72B AS PLANNER AND GROUNDER

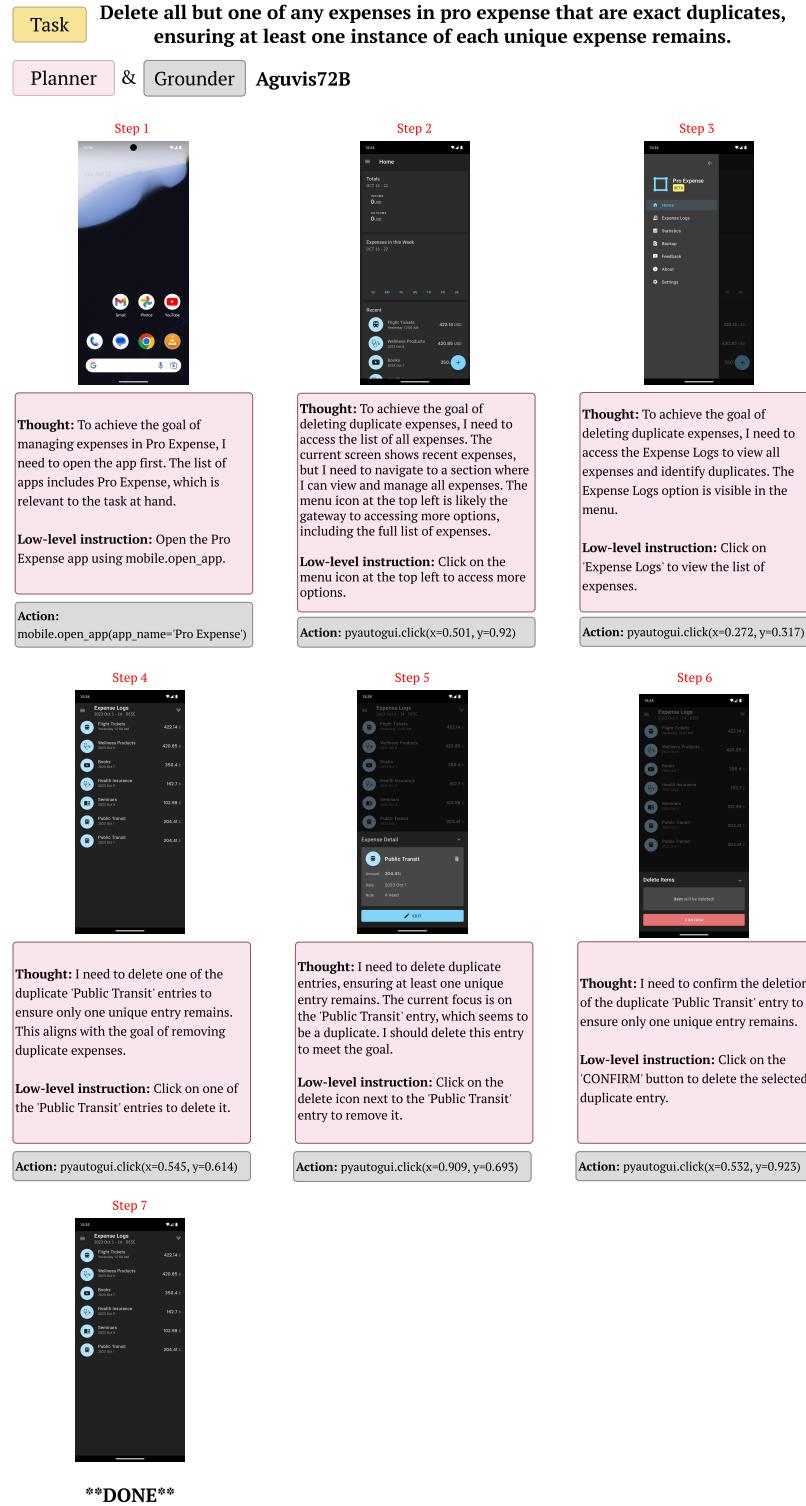


Figure 9: Example of AGUVIS-72B as planner and grounder executing AndroidWorld task. Due to space limitations, we present here the trajectory generated guided by Thought.

E.3.4 ANDROIDWORLD CASE: GPT-4O AS PLANNER AND AGUVIS-7B AS GROUNDER

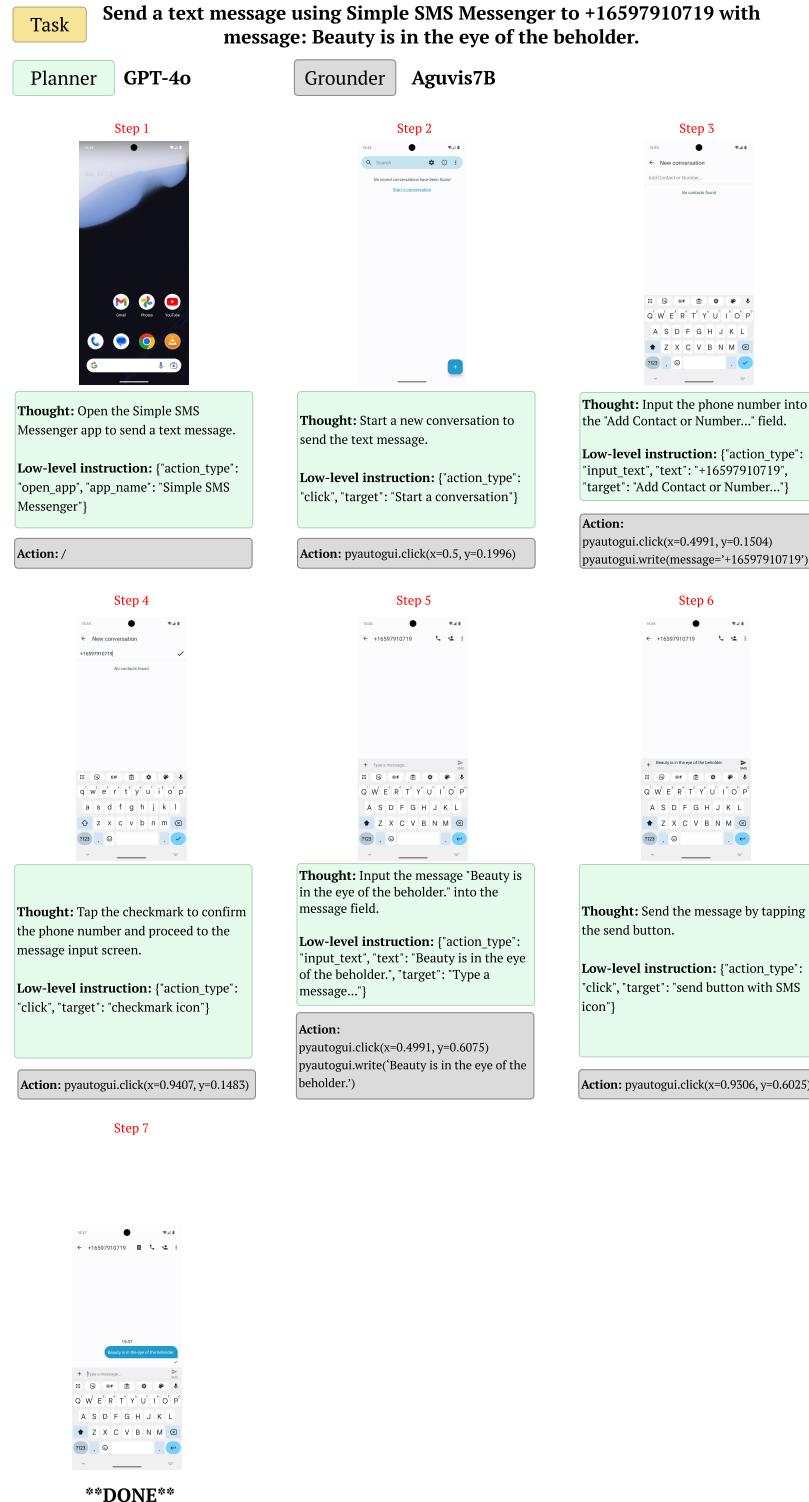
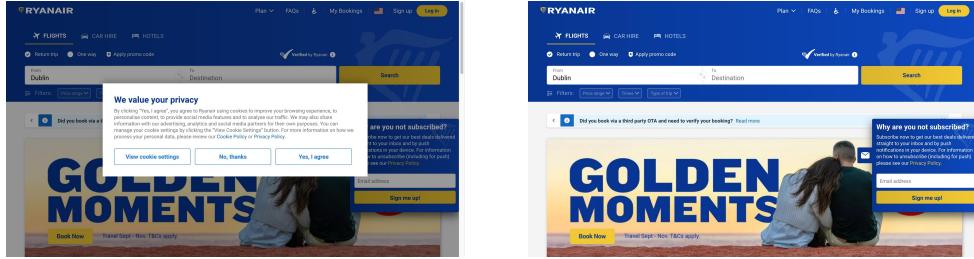


Figure 10: Example of GPT-4o as planner and AGUVIS-7B as grounder executing AndroidWorld task.

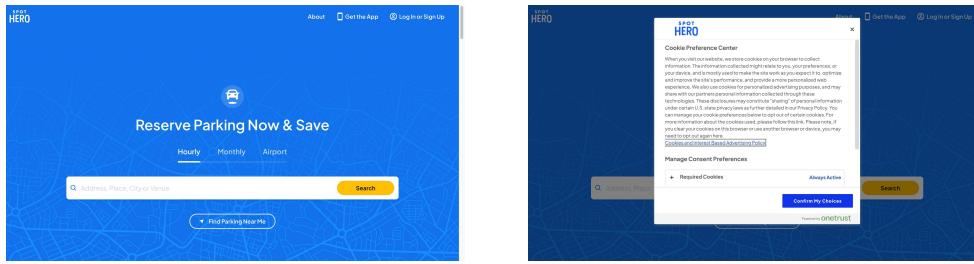
E.4 CASE OF AGUVIS GENERALIZATION IN REAL-WORLD SCENARIOS



Thought: To proceed with finding a flight, I need to dismiss the privacy pop-up to access the search functionality.

Low-level instruction: Click 'No, thanks' to dismiss the privacy pop-up.

Action: pyautogui.click(x=0.498, y=0.603)



Thought: To proceed with contacting support, I need to close this overlay to access the main page where the 'Contact Us' or 'Support' link is likely located.

Low-level instruction: Click the 'X' button to close the cookie consent overlay.

Action: pyautogui.click(x=0.704, y=0.093)

Figure 11: Case of AGUVIS generalization in real-world scenarios: closing cookie pop-ups, which is an out-of-domain situation.