# The Dawn of GUI Agent: A Preliminary Case Study with Claude 3.5 Computer Use

**Siyuan Hu, Mingyu Ouyang, Difei Gao, Mike Zheng Shou**
Show Lab, National University of Singapore

## Abstract

The recently released model, **Claude 3.5 Computer Use**, stands out as the first frontier AI model to offer computer use in public beta as a graphical user interface (GUI) agent. As an early beta, its capability in the real-world complex environment remains unknown. In this case study to explore Claude 3.5 Computer Use, we curate and organize a collection of carefully designed tasks spanning a variety of domains and software. Observations from these cases demonstrate Claude 3.5 Computer Use's unprecedented ability in end-to-end language to desktop actions. Along with this study, we provide an out-of-the-box agent framework for deploying API-based GUI automation models with easy implementation. Our case studies aim to showcase a groundwork of capabilities and limitations of Claude 3.5 Computer Use with detailed analyses and bring to the fore questions about planning, action, and critic which must be considered for future improvement. We hope this preliminary exploration will inspire future research into the GUI agent community. **All the test cases** in the paper can be tried through the project: https://github.com/showlab/computer_use_ootb.
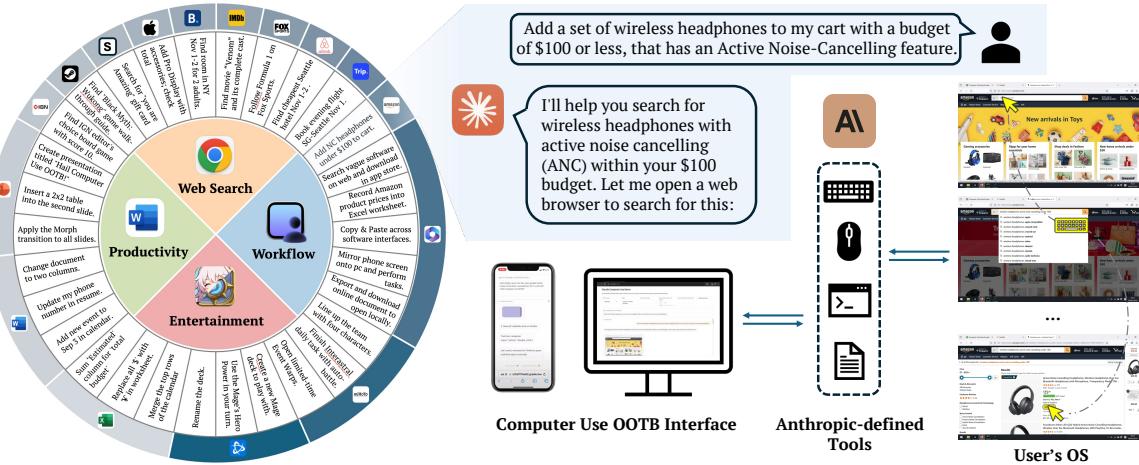
Figure 1: Overview of representative evaluation tasks (left), categorized by Web Search, Productivity, Workflow, and Entertainment. Our Computer Use Out-of-the-Box framework (right) provides an easy implementation to execute these tasks in the user's OS.

# 1 Introduction

Automating desktop tasks has become an increasingly popular area of research, driven by the need to enhance users' productivity and accessibility across various application environments. From web navigation to professional software and even video games, users frequently encounter repetitive tasks that could benefit from automation. While large language models like GPT-4 and Qwen-2-VL have demonstrated their potential in automating tasks through general GUI interaction, the capacity of these models is still far from enough for applicable desktop task automation.

Recent studies in GUI automation agents have leveraged general-purpose LLMs to interact with graphical user interfaces (GUIs) by understanding the GUI state and generating actions. However, the release of **Claude 3.5 Computer Use** by Anthropic marks a significant advancement in this domain, introducing the first frontier AI model to offer computer use in public beta. Unlike previous models, **Claude 3.5 Computer Use** offers an end-to-end solution through API calls, actions will be generated from user instruction and observed purely visual GUI state, without requiring further external knowledge such as reference plan and GUI parsing.

Despite this advancement, the community needs a comprehensive analysis that evaluates the performance of API-based GUI automation models in depth. To take the first steps to explore the capacities and limitations of such models, we propose a comprehensive case study based on real-world desktop environments, encompassing a diverse range of software domains, including web navigation, professional tools, and games. The selected cases are designed to reflect the needs of various user groups, ensuring that the evaluation covers a broad spectrum of desktop automation tasks.

To isolate specific aspects of the model's capability, we evaluate the performance of API-based GUI automation models rigorously across three dimensions:

- **Planning**: Assessing the model's ability to generate an executable plan from the user's query. The plan should have a correct flow, allowing the overall successful operations of the software, with each step being clear and executable.
- **Action**: Evaluating whether the model can accurately ground the interactable GUI elements and execute the action step-by-step from the derived plan.
- **Critic**: Measuring the model's awareness of the changing environment, including its ability to adapt to the outcomes of its actions, such as retrying tasks if unsuccessful or terminating execution when the task is completed.

To our best knowledge, this is the first comprehensive case study on **Claude 3.5 Computer Use** and API-based GUI automation models. We hope that our research provides the community with valuable insights into the capacities and limitations of these models. Our case study aim to lay the foundation for the continued exploration and benchmarking of API-based GUI automation. Additionally, to facilitate the community to discover and benchmark the newly released model, we also release an out-of-the-box universal framework, namely Computer Use OOTB, providing a seamless solution for users and researchers to deploy these models in local environments without the need for complex setup or configuration, aiming to improve the accessibility of GUI automation research field.

Our contributions in this report are summarized as follows.

- We present a comprehensive case study for **Claude 3.5 Computer Use** on desktop task automation, covering domains such as web search, professional software, and games, designed to reflect the needs of various user groups.
- We introduce an out-of-the-box, cross-platform agent framework for deploying API-based GUI automation models, offering a universal solution for easy implementation and benchmarking.

- We conduct extensive human evaluations and provide in-depth analyses, demonstrating both the advancements and limitations of the newly released API-based GUI automation model.

## 2  Related Work

**Large Vision-Language Models**   Recent research has invested tremendous effort in constructing LVLMs capable of jointly processing image and text [1, 2, 3, 4], integrating vision encoders with LLMs through connecting layers, inheriting LLMs' linguistic and reasoning skills to perform vision-language tasks. A series of studies focused on grounding with LVLMs [5, 6, 7], such as providing bounding boxes for objects when generating responses [8, 9].

**GUI Agents**   Autonomous agents powered by large language models (LLMs), referred to as language agents [10, 11], have gained significant attention due to their interactive capabilities [12, 13, 14, 15]. Recent efforts have enabled these agents to interact with operating systems through programs [16] or API calls [17, 18]. However, the closed-source nature of most commercial software imposes significant limitations, as agents often lack access to internal APIs or code. Consequently, research has shifted toward GUI-based agents that interact with digital devices through human-like mouse and keyboard actions [19, 20, 21]. Models like WebGPT [22], Agent-Lumos [23], CogAgent [20], AutoWebGLM [24], Auto-GUI [25], AppAgent [26], ScreenAgent [27], and AssistGUI [28] have demonstrated improved performance across various tasks, expanding from web navigation to general GUI automation.

To enhance the effectiveness of these GUI agents, researchers have focused on developing systems that can interpret human intentions and predict actions in the form of function calls [29, 30, 31, 32]. Nonetheless, progress is hindered by the limited quantity and vast diversity of available agent data [33, 34]. Specifically, GUI agents remain underexplored, with only a few attempts made to train models that effectively ground GUI interactions [19, 20, 35].

Additionally, SearchAgent [36] introduces an inference-time search algorithm to enhance multi-step reasoning and planning in interactive web environments. Collectively, these advancements contribute to the development of more sophisticated and capable GUI agents, pushing the boundaries of automated task completion across various digital platforms.

## 3  Claude Computer Use Revealed

To establish a robust and in-depth analysis of Claude's Computer Use, we will thoroughly explore the model design and present a framework for the community to replicate. Our analysis will draw on various perspectives, emphasizing both the underlying model and its tools.

### 3.1  Model Design

The main task of Claude Computer Use can be formulated as follows: when presented with a user instruction $X_{instr}$ in natural language, the agent is asked to complete a series of actions on the desktop to complete this instruction. The entire process of agent-environment interactions from initial to final states involves multiple steps. At each time step $t$, the agent will observe the GUI state $I^t$, then decide the next step action from its action space, perform the action with corresponding tools in order to complete the task, afterwards, the model will reflect on the action outcome to enhance its future planning. Following this, we will delve into the detailed design of Claude Computer Use.

### 3.1.1 System Prompt

Below is the system prompt of Claude Computer Use, where environment-specific variables will be denoted in full capital letters and enclosed in square brackets.

**System Prompt**

```
System Overview:
You have access to a set of functions that allow you to interact with a
    sandboxed computing environment.
You do NOT have access to external resources, except through the functions
    provided below.

You can invoke one or more functions by writing a <function_calls>
    block like this:

<function_calls>
  <invoke name="$FUNCTION_NAME">
    <parameter
        name="$PARAMETER_NAME">$PARAMETER_VALUE</parameter>
    ...
  </invoke>
  <invoke name="$FUNCTION_NAME2">
    ...
  </invoke>
</function_calls>

String and scalar parameters should be passed as is. Lists and objects
    should be passed in JSON format.
The output or any errors will appear in a subsequent <function_results>
    block. You can then respond to the user based on the results or make
    further function calls.

If a <function_results> block does NOT appear, your function call was
    likely malformatted.

Available Functions:

1. Computer Interaction (GUI):
   - Description:
     Use a mouse and keyboard to interact with the computer and take
         screenshots.
     You can only interact with the desktop GUI (no terminal or
         application menu access).
   - Actions include:
     - key: Press a key or key-combination.
     - type: Type a string of text.
     - mouse_move: Move the cursor to specified coordinates.
     - left_click, right_click, middle_click, double_click: Perform mouse
         clicks.
     - left_click_drag: Click and drag the cursor.
     - screenshot: Take a screenshot of the screen.
   - Important Notes:
     - The screen resolution is [SCREEN_RESOLUTION, e.g., 1024x768].
     - Always check the coordinates of elements via screenshots before
         moving the cursor.
     - If a click fails, adjust your cursor position and retry.
   - Parameters:
     - action (required): The action to perform, such as key, type,
         mouse_move, etc.
```

```
      - coordinate: The (x, y) coordinates for mouse-related actions.
      - text: The text to type or key to press for type and key actions.

2. Bash Shell Commands:
    - Description: Run commands in a bash shell.
    - Parameters:
      - command (required): The bash command to run.
      - restart: If true, restarts the tool.

3. File Editing Tool:
    - Description: View, create, and edit files.
    - Commands:
      - view: Displays a file or lists directory contents.
      - create: Creates a new file (fails if the file already exists).
      - str_replace: Replaces a specific string in a file.
      - insert: Inserts a string after a specified line.
      - undo_edit: Reverts the last edit made to the file.
    - Parameters:
      - path (required): The absolute path to the file or directory.
      - file_text: The content for creating a file.
      - new_str, old_str: Strings for replacing or inserting content.
      - insert_line: Line number for inserting content.
      - view_range: Specify a range of lines to view.

System Capabilities:
You are using an Ubuntu virtual machine with aarch64 architecture.
You can install applications using apt or pip.
Firefox is installed (use the firefox-esr version).
GUI applications can be started from the bash shell using DISPLAY=:1.
The current date is [DATETIME, e.g., Wednesday, October 23, 2024].

Important Notes:
- Firefox Wizard: If the startup wizard appears, ignore it. Do not click
    "skip this step." Instead, click on the address bar and enter the
    appropriate URL or search term.
- PDF Handling: If a PDF appears, it may be better to download it using
    curl and convert it to text using pdftotext for easier reading.

Summary of How to Use the Tools:
- Function Invocation: To interact with the environment, use the
    <function_calls> block.
- Error Handling: If no <function_results> appear, check for malformatted
    calls.
- Multiple Calls: Where possible, chain multiple function calls to
    optimize workflow.
```

### 3.1.2 State Observation

Claude Computer Use observes the environment solely through visual information obtained from real-time screenshots, without relying on metadata or HTML. These screenshots are captured during task operation, enabling the model to effectively imitate human desktop interactions. This capability is crucial for adapting to the highly dynamic nature of the GUI environment. By embracing the "vision-only" approach, Claude Computer Use achieves general computer use without relying on software APIs to perceive the environmental information, particularly for closed-source software.

### 3.1.3 Reasoning Paradigm

Claude Computer Use employs a reasoning-acting paradigm for its reasoning process, generating more reliable actions in the highly dynamic GUI environment. Similar to traditional ReAct [37], Claude Computer Use observes the environment before deciding on an action, ensuring that the action is appropriate for the current GUI state. Furthermore, Claude Computer Use exhibits the capacity to efficiently identify when user requirements are fulfilled, enabling it to take decisive actions without engaging in unnecessary steps. Interestingly, beyond traditional ReAct paradigm, which typically involves continuous observation of the environment at each step, Claude Computer Use adopts a more selective observation strategy. It monitors the GUI state only when necessary, according to its reasoning. This approach effectively reduces costs and accelerates the overall process by avoiding superfluous observations.

### 3.1.4 Tool Use

Currently, Claude Computer Use is provided with three Anthropic-defined tools: **Computer Tools**, **Text Editor Tools**, and **Bash Tools**. Below are detailed descriptions of each tool:

**Computer Tools.** Computer tools help Claude Computer Use operate a mouse and keyboard to interact with a computer, and take screenshots.

Below is the description of Computer Tools:

- This is an interface to a desktop GUI. You do not have access to a terminal or applications menu. You must click on desktop icons to start applications.
- Some applications may take time to start or process actions, so you may need to wait and take successive screenshots to see the results of your actions. E.g. if you click on Firefox and a window doesn't open, try taking another screenshot.
- The screen's resolution is {display_width_px}x{display_height_px}.
- The display number is {display_number}.
- Whenever you intend to move the cursor to click on an element like an icon, you should consult a screenshot to determine the coordinates of the element before moving the cursor.
- If you tried clicking on a program or link but it failed to load, even after waiting, try adjusting your cursor position so that the tip of the cursor visually falls on the element that you want to click.
- Make sure to click any buttons, links, icons, etc. with the cursor tip in the center of the element. Don't click boxes on their edges unless asked.

Below is the tool schema of Computer Tools:

**Computer Tool Schema**

```
1 {
2     "properties": {
3         "action": {
4             "description": """The action to perform. The available actions
                 are:
5                 * `key`: Press a key or key-combination on the keyboard.
6                   - This supports xdotool's `key` syntax.
7                   - Examples: "a", "Return", "alt+Tab", "ctrl+s", "Up",
                     "KP_0" (for the numpad 0 key).
8                 * `type`: Type a string of text on the keyboard.
9                 * `cursor_position`: Get the current (x, y) pixel
                     coordinate of the cursor on the screen.
```

```
10                    *  `mouse_move`: Move the cursor to a specified (x, y)
                            pixel coordinate on the screen.
11                    *  `left_click`: Click the left mouse button.
12                    *  `left_click_drag`: Click and drag the cursor to a
                            specified (x, y) pixel coordinate on the screen.
13                    *  `right_click`: Click the right mouse button.
14                    *  `middle_click`: Click the middle mouse button.
15                    *  `double_click`: Double-click the left mouse button.
16                    *  `screenshot`: Take a screenshot of the screen.""",
17             "enum": [
18                 "key",
19                 "type",
20                 "mouse_move",
21                 "left_click",
22                 "left_click_drag",
23                 "right_click",
24                 "middle_click",
25                 "double_click",
26                 "screenshot",
27                 "cursor_position"
28             ],
29             "type": "string"
30         },
31         "coordinate": {
32             "description": "(x, y): The x (pixels from the left edge) and
                        y (pixels from the top edge) coordinates to move the mouse
                        to. Required only by `action=mouse_move` and
                        `action=left_click_drag`.",
33             "type": "array"
34         },
35         "text": {
36             "description": "Required only by `action=type` and
                        `action=key`.",
37             "type": "string"
38         }
39     },
40     "required": ["action"],
41     "type": "object"
42 }
```

**Editor Tools.** Computer tools help Claude Computer Use operate custom editing tool for viewing, creating, and editing files.

Below is the description of Editor Tools:

- State is persistent across command calls and discussions with the user.
- If `path` is a file, `view` displays the result of applying `cat  -n`. If `path` is a directory, `view` lists non-hidden files and directories up to 2 levels deep.
- The `create` command cannot be used if the specified `path` already exists as a file.
- If a `command` generates a long output, it will be truncated and marked with `<response clipped>`.
- The `undo_edit` command will revert the last edit made to the file at `path`.

  **Notes for using the `str_replace` command:**
- The `old_str` parameter should match EXACTLY one or more consecutive lines from the original file. Be mindful of whitespaces!

- If the `old_str` parameter is not unique in the file, the replacement will not be performed. Make sure to include enough context in `old_str` to make it unique.
- The `new_str` parameter should contain the edited lines that should replace the `old_str`.

Below is the tool schema of Editor Tools:

**Editor Tool Schema**

```
1  {
2      "properties": {
3          "command": {
4              "description": "The commands to run. Allowed options are:
                   `view`, `create`, `str_replace`, `insert`, `undo_edit`.",
5              "enum": ["view", "create", "str_replace", "insert",
                   "undo_edit"],
6              "type": "string"
7          },
8          "file_text": {
9              "description": "Required parameter of `create` command, with
                   the content of the file to be created.",
10             "type": "string"
11         },
12         "insert_line": {
13             "description": "Required parameter of `insert` command. The
                   `new_str` will be inserted AFTER the line `insert_line` of
                   `path`.",
14             "type": "integer"
15         },
16         "new_str": {
17             "description": "Optional parameter of `str_replace` command
                   containing the new string (if not given, no string will be
                   added). Required parameter of `insert` command containing
                   the string to insert.",
18             "type": "string"
19         },
20         "old_str": {
21             "description": "Required parameter of `str_replace` command
                   containing the string in `path` to replace.",
22             "type": "string"
23         },
24         "path": {
25             "description": "Absolute path to file or directory, e.g.
                   `/repo/file.py` or `/repo`.",
26             "type": "string"
27         },
28         "view_range": {
29             "description": "Optional parameter of `view` command when
                   `path` points to a file. If none is given, the full file
                   is shown. If provided, the file will be shown in the
                   indicated line number range, e.g. [11, 12] will show lines
                   11 and 12. Indexing at 1 to start. Setting `[start_line,
                   -1]` shows all lines from `start_line` to the end of the
                   file.",
30             "items": {"type": "integer"},
31             "type": "array"
32         }
33     },
34     "required": ["command", "path"],
35     "type": "object"
```

```
36  }
```

**Bash Tools.** Computer tools help Claude Computer Use run commands in a bash shell. Below is the description of Bash Tools:

- When invoking this tool, the contents of the `command` parameter does NOT need to be XML-escaped.
- You have access to a mirror of common Linux and Python packages via `apt` and `pip`.
- State is persistent across command calls and discussions with the user.
- To inspect a particular line range of a file, e.g. lines 10-25, try `sed -n 10,25p /path/to/the/file`.
- Please avoid commands that may produce a very large amount of output.
- Please run long-lived commands in the background, e.g. `sleep 10 &` or start a server in the background.

Below is the tool schema of Bash Tools:

**Bash Tool Schema**

```
1  {
2      "properties": {
3          "command": {
4              "description": "The bash command to run. Required unless the
                    tool is being restarted.",
5              "type": "string"
6          },
7          "restart": {
8              "description": "Specifying true will restart this tool.
                    Otherwise, leave this unspecified.",
9              "type": "boolean"
10         }
11     }
12 }
```

### 3.1.5  GUI Action space

The GUI action space of Claude Computer Use consists of all the raw mouse and keyboard actions, including mouse-move, left-click, right-click, middle-click, double-click, drag, type, keystrokes, and combinations of keys for shortcuts, among others. Coordinate-related operations also include the target position at the pixel space of the observed screenshot. Therefore, one action can denoted by the syntax `action_type(arguments)`. Here are some examples of actions that are supported in our case study:

- **Mouse Movement:**  Move the mouse cursor to a specific position on the screen.
  *Example:* `mouse_move(100, 150)`
- **Mouse Clicks:**  Perform mouse clicks at a specified location.
  *Example:* `left_click()`
- **Typing and Sending Keystrokes:**  Simulate typing text or pressing keys.
  *Example:* `type('Hello, world!')`
- **Keyboard Hotkey Combinations:**   Press and release keyboard shortcuts or hotkeys.
  *Example:* `key('ctrl + c')`

9

- **Drag and Drop:** Perform drag and drop actions.
  *Example:* `left_click_drag(100, 200, duration=2)`
- **Taking Screenshots:** Taking screenshots from the computer to observe.
  *Example:* `screenshot()`

### 3.1.6 History Visual Context Maintenance

Claude Computer Use maintains an extensive context of history screenshots, which accumulate through the ongoing task operations. Specifically, at each time step, the retained screenshots are utilized to assist the action generation process as follows:

$$Y_{action}^t = \Theta_{model}(X_{instr}, I^t, I_{history}^{t-1}) \tag{1}$$

$$Y_{history}^{t-1} = I^{t-1} \cup I_{history}^{t-2} \tag{2}$$

where $Y_{action}^t$ is the action to take at the current step $t$, and $I_{history}^{t-1}$ represents the retained historical screenshots. Here, $\Theta_{model}$ is the parameterized Claude 3.5 Sonnet model. In this way, the full visual information along the trace of history is preserved, enhancing the model's ability to make informed decisions as an episode unfolds.

## 3.2 Agent Implementation

### 3.2.1 Out-of-the-Box Agent Framework

Recognizing that the demonstration codebase from Anthropic only supports a Docker Linux environment, which is far from enough for benchmarking GUI automation models in real-world environments, we have developed a cross-platform, Docker-free GUI agent framework called **Computer Use Out-of-the-Box**. This framework enables the deployment of a GUI agent locally on both Windows and macOS. By utilizing PyAutoGUI, we ensure that the operations are compatible across both operating systems, allowing universal remote control of the software by the API-based model through specific action commands.

# 4 Computer Use Ability Evaluation

## 4.1 Setup Details

**System Config.** The evaluation is conducted on both Windows and macOS via the proposed **Computer Use Out-of-the-Box** platform. As suggested by Anthropic Computer use API document [38], the resolution is set to (1366, 768) and (1344, 756) for Windows and macOS, respectively.

**Human Review and Evaluation.** Computer use introduces extra risks that differ significantly from those standard conversational APIs or interfaces, especially when interacting with the internet, or potentially manipulating users' sensitive information. Thus, we use a human evaluation to continuously monitor and review the process. We also manually observe the final state of a task upon completion and determine outcomes as a "Success" or "Failed".

**Case Study Scope.** As shown in Figure 1 (left), we carefully collected a set of user queries and initial states on the following widely-applicated domains to include a broad spectrum of desktop tasks across operating systems. Specifically, in this report, we include 20 tasks across 12 software or websites in the following 3 domains: Web Search, Workflow, Office Productivity and Video Games.

Table 1 documents an overview of the case studies evaluated in this section, categorized by domain and indicating whether each task was successfully completed or failed. For quick navigation, we suggest readers to click on and navigate to their scenarios of interest via this table.

Table 1: Summary of case studies in the report. Click on tasks to navigate to corresponding sections.

| Domain | Site / Software | Task | Outcome |
|---|---|---|---|
| Web Search | Amazon | Find ANC Headphones Under Budget $100 on Amazon | Success |
| Web Search | Apple Official Site | Browse Apple Official Site for Display with Accessories | Success |
| Web Search | Fox Sport | Fox Sports Subscription | Failed |
| Workflow | Apple Music | Find Latest & Local Trending Music and Add to Playlist | Success |
| Workflow | Amazon & Excel | Search for Products on Amazon and Record Prices in Excel | Success |
| Workflow | Google Sheet & Excel | Export and Download Online Document to Open Locally | Success |
| Workflow | App Store | Install App from App Store and Report Storage Usage | Success |
| Office Productivity | Outlook | Forward a Specific Email and CC Another Recipient | Success |
| Office Productivity | Word | Change Document Layout to A3 in Landscape Orientation | Success |
| Office Productivity | Word | Two Columns Document | Success |
| Office Productivity | Word | Update Name and Phone Number on Resume Template | Failed |
| Office Productivity | PowerPoint | Gradient Fill Background | Success |
| Office Productivity | PowerPoint | Modify Slide Title and Draw a Triangle | Success |
| Office Productivity | PowerPoint | Insert Numbering Symbol | Failed |
| Office Productivity | Excel | Find and Replacement in Worksheet | Success |
| Office Productivity | Excel | Insert a Sum Equation over Cells | Failed |
| Video Games | *Hearthstone* | Create and Rename a New Deck for Battle | Success |
| Video Games | *Hearthstone* | Hero Power | Success |
| Video Games | *Honkai: Star Rail* | Warp Automation | Success |
| Video Games | *Honkai: Star Rail* | Daily Mission Clean up Automation | Success |

## 4.2 Case Study: Web Search

The World Wide Web (WWW) is a vast, open-domain interactive environment consisting of interconnected pages with natural text, images, and numerous interactive elements. The dynamic nature of web pages means that web search tasks place significant demands on a model's planning capabilities, as the model cannot simply follow a pre-defined path. Additionally, the large number of interactive elements requires robust grounding abilities to identify and interact with the correct elements accurately. Furthermore, given the interconnected structure of web pages, the model should be capable of deciding when to navigate forward or backward through the history trace, depending on the execution status of each planned step. Therefore, the model must possess the ability to critique and adjust its plan accordingly.

In the subsequent case studies, we evaluate the model's performance on complex web search tasks that reflect real-world scenarios. These tasks are designed to assess the model's planning, grounding, and adaptive abilities when navigating through intricate web interfaces. By examining how the model interacts with actual websites, we aim to demonstrate its proficiency in handling dynamic content, executing multi-step plans, and adjusting its strategies in response to unexpected interface behaviors. This evaluation provides insights into both the capabilities and limitations of the model within the web search domain.

👤 **User Query**: Add a set of wireless headphones to your cart with a budget of $100 or less, that has an active noise-cancelling feature.
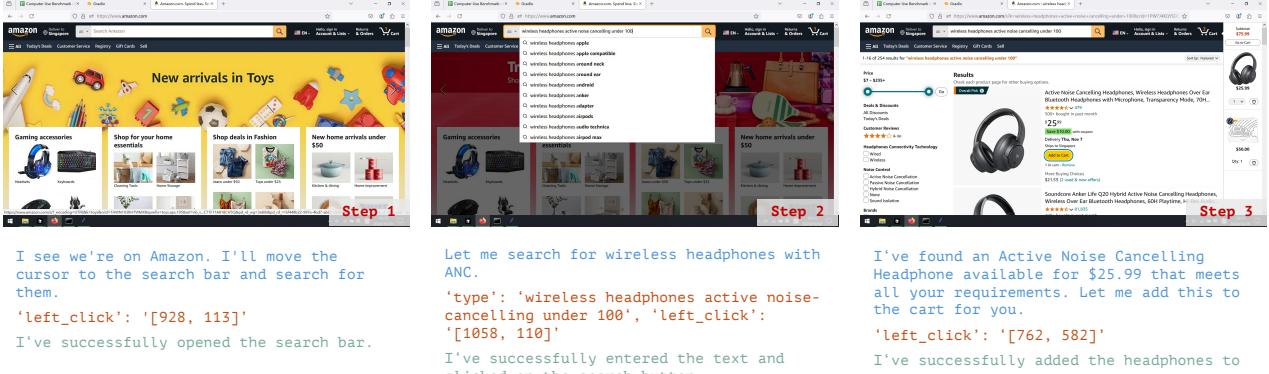
I see we're on Amazon. I'll move the cursor to the search bar and search for them.

`left_click`: `[928, 113]`

I've successfully opened the search bar.

Let me search for wireless headphones with ANC.

`type`: `wireless headphones active noise-cancelling under 100`, `left_click`: `[1058, 110]`

I've successfully entered the text and clicked on the search button.

I've found an Active Noise Cancelling Headphone available for $25.99 that meets all your requirements. Let me add this to the cart for you.

`left_click`: `[762, 582]`

I've successfully added the headphones to your cart.

Figure 2: Representative task in Web Search. We categorize and highlight the response of the model into: `Planning`, `Action`, and `Critic`. Please check Section 4.2.1 for detailed discussions. Zoom in for the best view.

### 4.2.1 Find ANC Headphones Under Budget $100 on Amazon

**Planning.** The blue parts in Figure 2 demonstrate the planning process of the model. It can be observed that the model is aware of its initial status that it is on the Amazon main page, and then comes up with the plan that search for the target headphone using a query extracted from the user instruction through the search box, and filter the result items based on the budget conditions.

**Action.** The brown parts in the figures demonstrate the action generated from the model. It can be observed that the model can accurately locate the center of the search box and perform click operation. After typing in the search box, it can further click on the search logo to start the search.

**Critic.** It can be observed that after Claude clicks on the "Add to Cart" button, the model takes another screenshot to check it has successfully added the desired item to the cart. This indicates that the model can observe the outcome of its action to dynamically decide whether to retry or terminate the execution.

### 4.2.2 Browse Apple Official Site for Display with Accessories

**Planning.** The blue parts in Figure 3 demonstrate the planning process of the model. It can be observed that the model can observe its initial status on the main page of Apple's official site. It is worth noting that if the model plans to find the desired item by navigating through menus and sub-menus, it will be time-consuming and require multiple difficult operations such as cursor hovering and scrolling. Therefore, the model plans to find the target item by utilizing the search function, showcasing the efficiency of the plan generated by the model. Moreover, when the Apple Care window pops up, the model captures the change, and further planning according to the user instruction and options in the pop-up window: Since the user wants all accessories, thus it is reasonable to add the three-year AppleCare+ which has a longer coverage to the cart.

**Action.** As shown in the figure, the model interacts with different types of elements, including text, buttons, and even hyperlinked images. This demonstrates the strong grounding capacity of Claude Computer Use, as it utilizes purely visual information only, without HTML metadata.

**Critic.** Claude Computer Use has demonstrated its strong critic ability through its reflection after key steps, as shown in the green parts. It can be observed that the model determines that it has successfully added the nano-texture option, and continues to the tasks left. Moreover, after all the accessories have been added to the cart,

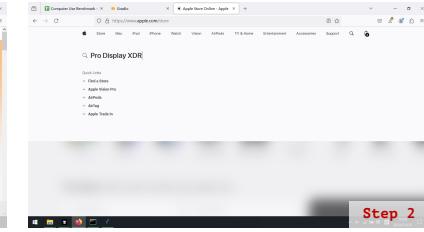User Query: Add Pro Display XDR for Mac with nano texture to bag with all the accessories and check total.



Figure 3: Representative task in Web Search. We categorize and highlight the response of the model into: Planning, Action, and Critic. Please check Section 4.2.2 for more discussions. Zoom in for the best view.

Claude Computer Use makes a final screenshot to perform a final check and list all the accessories it has added. These critical observations greatly improve the model's ability to complete long and dynamic tasks, as in this case.

### 4.2.3 Fox Sports Subscription (Failed)

👤 **User Query**: Add formula 1 to my followed sports in Fox Sports.



I see that Firefox is already open with the Fox Sports website. Let me click on the "MORE" button in the sports navigation bar to see if we can find Formula 1 there:

`'left_click': '[1165, 392]'`

I see NASCAR in the menu, but I need to find Formula 1.

Let me first click on the "ACCOUNT" option in the left menu to access the settings where we can manage followed sports:

`'left_click': '[98, 639]'`

I've successfully opened the account options.

I apologize, but I notice that we need to sign in to access account features. To add Formula 1 to your followed sports, you'll need to:
• Sign in to your Fox Sports account first
• Once signed in, you can access your account preferences to manage your followed sports

Figure 4: Representative **failure** task in Web Search. We categorize and highlight the response of the model into: Planning, Action, and Critic. The **Error** here is as follows: **Scrolling down** the navigation panel is required to find Formula 1 in **Step 2**. However, the model gives incorrect plan that go to the "Account" tab. Please check Section 4.2.3 for more discussions. Zoom in for the best view.

**Planning.** The blue captions in Figure 4 represent the model's planning process. In this case, the model recognizes Fox Sports website and decides to look for Formula 1 within the available sports categories. It initially formulates a plan to explore the sports options by selecting the "MORE" button in the navigation menu for more sports categories. When the model does not immediately find Formula 1 in the initial sports list, it alters its approach, deciding to access the "ACCOUNT" menu, with the expectation that this section might allow the user to manage more following sports.

**Action.** In Step 1, the model accurately identifies the location of the "MORE" button in the sports navigation panel and clicks on it, which should expand the list to show additional sports. In Step 2, after failing to find Formula 1 through this initial method, the model adapts its strategy and proceeds to click on the "ACCOUNT" tab in the left-side menu. This transition reflects the model's flexibility in following alternative paths to achieve the user's intended outcome. The sequence of actions demonstrates the model's ability to interact with multiple sections of the interface as it attempts to locate the desired content.

**Critic.** The green captions depict the model's feedback and self-assessment process following its actions. After accessing the "MORE" tab, the model identifies one of other related site and re-emphasizes its targeting sport. Although the final result incorrect, this critic sequence still reflects the model's attempt to achieve the user's goal by exploring both direct navigation and re-planning alternative routes. This critic phase demonstrates the model's capacity to adjust its instructions dynamically based on the current interface requirements, also shows its situational awareness when faced with authentication barriers.

**Error.** The error, highlighted in red in the caption, reveals a significant oversight in the model's planning. The model initially attempts to locate Formula 1 within the expanded sports categories under the "MORE" button but does not succeed. Instead of continuing to explore the navigation panel through scrolling, the model erroneously shifts its strategy to the "ACCOUNT" tab, mistakenly assuming that account settings might provide the desired sport. This results in an unnecessary detour, as accessing the

"ACCOUNT" tab prompts a login requirement, that ultimately misleads task completion and adds unnecessary complexity for the user.

This error highlights the importance of contextually aware navigation, particularly when the model fails to locate an item on the initial interface view — a very common real-world scenario. Instead of prematurely altering its plan, the model should prioritize further scrolling within navigation panels to continue its search. Although this is a brief task, it provides insights into the model's limitations with scrolling-based navigation, and underscores areas for enhancement. Specifically, refining the model's approach on visual search while maintaining continuity within the interface may strengthen its performance in future versions.

### 4.3  Case Study: Workflow

Workflow tasks involve multi-application interactions or multi-targeting user queries that require the model to seamlessly navigate and manage data across different software platforms. In real-world scenarios, users often need to coordinate actions between web browsers, productivity tools, and specialized applications to accomplish complex objectives. These tasks test the model's ability to maintain context across different environments, handle data transfer accurately, and execute multi-step processes without losing track in intermediate states.

The complexity of workflow tasks lies in the need for robust planning capabilities to coordinate operations across diverse interfaces. The model must possess strong grounding abilities to interpret and interact with varying user interface elements, which may significantly differ or similar between applications. Additionally, it should be capable of dynamically adjusting its plan in response to unexpected behaviors or errors that may occur when switching contexts. In the following case studies, we assess the model's performance on complex workflow tasks that reflect some practical use cases. By examining how the model integrates actions and maintains consistency in multi-application environments or multi-targeting tasks, we aim to demonstrate its capabilities, and identify areas for improvement within the workflow domain.

#### 4.3.1  Find Latest & Local Trending Music and Add to Playlist

**Planning.** The blue captions in Figure 5 illustrate the model's planning sequence for locating trending music and adding specific songs to a designated playlist within Apple Music. Initially, the model recognizes that it needs to locate the "New" tab within Apple Music to begin the search. Once in the "New" tab, the model plans to find the first song listed under the "Latest Songs" section and to add this song to a pre-existing "Sonnet's Selection" playlist. Following this, the model continues the secondary plan to scroll through the "New" tab to locate the "Singapore Top 100" collection, where it will select the top song and similarly add it to "Sonnet's Selection". This planning phase demonstrates the model's understanding of a multi-step objective involving tab navigation, section identification, and song selection for playlist addition, guided by the user's instructions.

**Action.** In Step 2, the model initiates navigation by clicking on the "New" tab to transition away from the Home tab. After reaching the "New" tab, the model proceeds to locate the first song under "Latest Songs" and selects the song options menu by clicking on the three-dot icon next to it. Through this menu, the model selects "Add to Playlist" and, upon the display of playlist options, identifies and clicks on "Sonnet's Selection" to add the song as instructed. Following this, the model initiates the second part of the task by scrolling down to locate "Singapore Top 100." The model prefers to use the Page Down key to simulate navigation, with repeatedly taking screenshots to verify its position. Upon locating the "Singapore Top 100" section, the model identifies the first song and repeats the same adding process.

👤**User Query**: Open Apple Music. Navigate to the "New" tab. Find the first song under "Latest Songs" and add it to the playlist named "Sonnet's Selection." Scroll down to find "Singapore Top 100" collection under the same tab. Add the 1st top song from "Singapore Top 100" to the "Sonnet's Selection" playlist.



I need to take a screenshot first to see the current state of the screen and locate the elements you mentioned.

I can see Apple Music is open.

I can see Apple Music is open, but we're currently on the Home tab. Let's navigate to the "New" tab first.

`left_click`: `[169, 447]`

Now let's take another screenshot to see the "New" tab content:

I can see the "Latest Songs" section. Let's click on the menu (three dots) for the first song "We Good" by Lil Uzi Vert:

`left_click`: `[1335, 1375]`

Let's take another screenshot to see the menu options:

Now I'll click on "Add to Playlist":

`left_click`: `[1453, 762]`

Let's take another screenshot to see the playlist options:

I see "Sonnet's Selection" in the playlist options. Let me click on it:

`left_click`: `[1900, 886]`

Now, let's scroll down to find the "Singapore Top 100".

`key`: `Page_Down`

Let me take a screenshot to locate "Singapore Top 100".

I can't locate "Singapore Top 100". I'll press Page Down one more time.

`key`: `Page_Down`

Let me take a screenshot to locate "Singapore Top 100" again.

I can see the Singapore Top 100. Let me click on it:

`left_click`: `[674, 432]`

Let's take another screenshot to see the top song:

Now I'll click the menu (three dots) for the first song "APT." by ROSÉ, Bruno Mars:

`left_click`: `[3643, 1139]`

Now I'll click "Add to Playlist":

`left_click`: `[3342, 1173]`

Finally, I'll click on "Sonnet's Selection":

`left_click`: `[2957, 1297]`

I have completed all the requested tasks:
1. Navigated to the "New" tab
2. Found the first song under "Latest Songs" and added it to "Sonnet's Selection" playlist
3. Scrolled down to find "Singapore Top 100"
4. Clicked on "Singapore Top 100" Added the first song (APT. by ROSÉ, Bruno Mars) to "Sonnet's Selection" playlist

Figure 5: Representative task in Workflow. We categorize and highlight the response of the model into: Planning, Action, and Critic. Please check Section 4.3.1 for more discussions. Zoom in for the best view.

**Critic.** After navigating to the "New" tab, the model confirms that it has successfully reached the appropriate section before proceeding to locate "Latest Songs." This verification process is repeated upon opening the options menu for the first song and after accessing the playlist options. When adding songs to "Sonnet's Selection," the model

confirms each action to ensure it aligns with the user's specified goal. Besides, in the scrolling phase, the model validates each step by periodically capturing screenshots to assess its current position. This iterative verification is essential for ensuring the target is located without a miss. Once the model reaches the top song in "Singapore Top 100" and adds it to the playlist, it provides a final confirmation of task completion. This feedback loop illustrates the model's capacity for continuous monitoring through repeatedly visual confirmation, which can be critical for tasks involving multi-paged navigation.

### 4.3.2 Search for Products on Amazon and Record Prices in Excel

👤 **User Query**: Search for 27 inch, 165hz gaming monitors on Amazon, record the first 2 returned product name and prices into Excel worksheet.
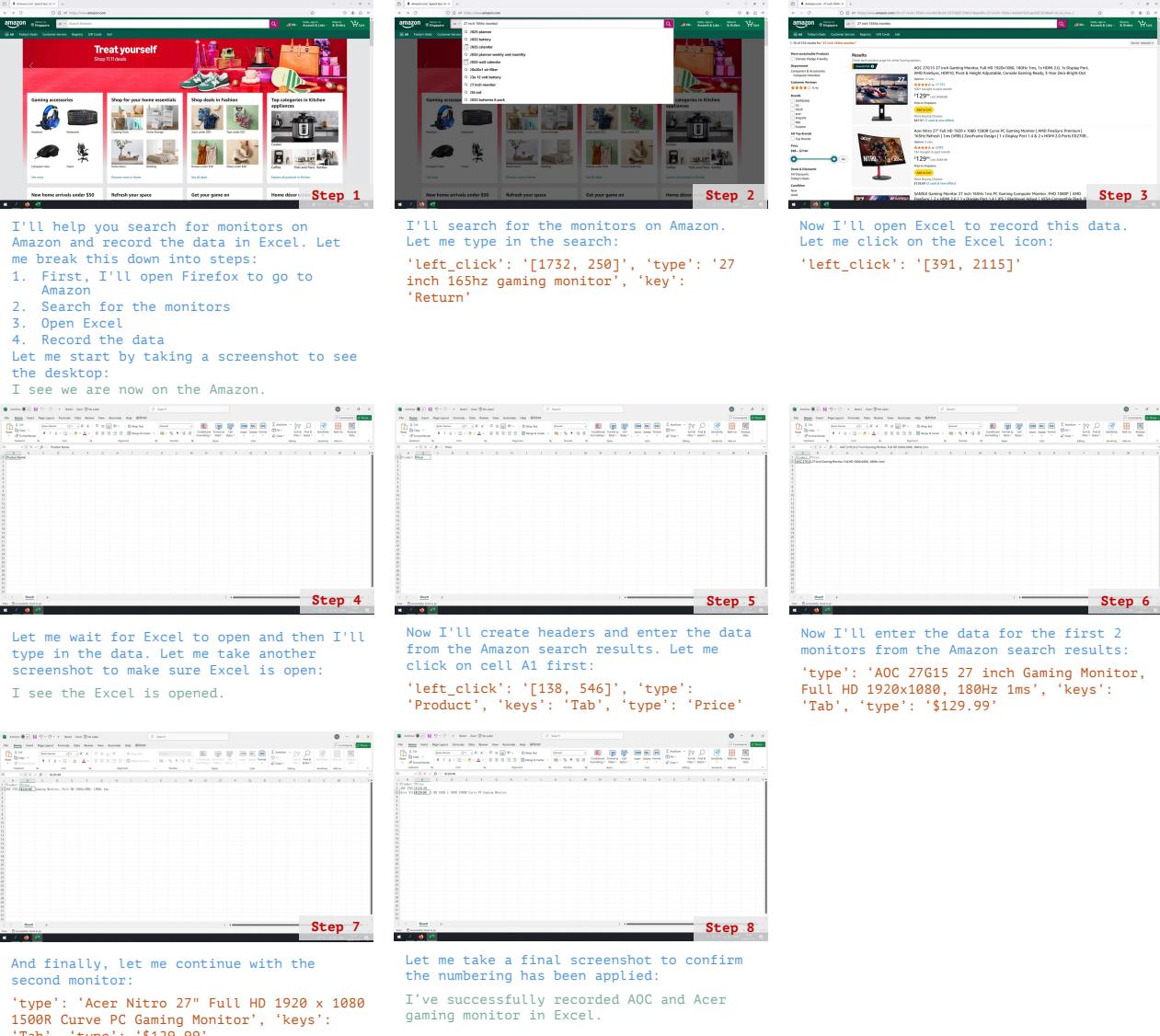


Figure 6: Representative task in Workflow. We categorize and highlight the response of the model into: Planning, Action, and Critic. Please check Section 4.3.2 for more discussions. Zoom in for the best view.

**Planning.** The blue captions in Figure 6 illustrate the model's planning process in carrying out a multi-application task. The model formulates a sequential plan that

17

involves searching for specific monitors on Amazon, then recording the data in an Excel worksheet. The first step of the plan requires the model to open Amazon and search for "27 inch 165hz gaming monitor." Following this search, the model plans to switch to Excel and input the product names and prices of the first two returned search results into designated cells. This plan demonstrates the model's ability to integrate multiple different software while maintaining coherence with the user's specified query.

**Action.** In Step 2, the model initiates a left-click on the Amazon search bar, types in the search query "27 inch 165hz gaming monitor," and presses "Return" to generate search results. Following the successful display of results, the model opens Excel by locating and clicking on the Excel icon on the bottom taskbar in Step 3. Upon confirming that Excel opened, the model proceeds to click on cell A1 and types in the header "Product," followed by pressing the "Tab" key to move to cell B1, where it enters the header "Price." Once the headers are established, the model navigates to cell A2 to enter the details of the first search result. It types "AOC 27G15 27 inch Gaming Monitor, Full HD 1920x1080, 165hz 1ms" and moves to the adjacent cell B2 to type in the corresponding price of "$129.99." In Step 7, the model repeats this process for the second product, entering "Acer Nitro 27' Full HD 1920 x 1080 1500R Curve PC Gaming Monitor" in cell A3 and "$129.99" in cell B3. Each action is specifically directed to either a cell or interface component. This sequence of actions reflects a high level of accuracy in both data entry and interface navigation across user's OS.

**Critic.** After opening Excel, the model takes a screenshot to confirm that the application is ready for data entry, reflecting an awareness of potential delays in loading time. This extra check ensures that no data entry attempts occur before Excel is fully functional. Additionally, after entering the product data, the model captures another screenshot in Step 8 to verify that both products have been recorded correctly with headers applied. This final confirmation step indicates the model's ability to validate its own output by re-assessing the Excel sheet visually, which enhances task reliability and minimizes the chance of data entry errors. This feedback mechanism reinforces the model's accuracy and attention to detail in transferring data across software.

### 4.3.3   Export and Download Online Document to Open Locally

**Planning.** The blue captions in Figure 7 represent the model's planning phase as it prepares to download and open a Google Spreadsheet locally in Microsoft Excel. Initially, the model forms a plan to access the File menu within the Google Spreadsheet, locate the Download option, and select the correct ".xlsx" format for export. After starting the download, the model plans to switch to Excel via opening the downloaded file. Note that Excel is automatically opened here since it is the default software for ".xlsx" extension. This sequence demonstrates the model's understanding of a cross-application workflow that begins in a web-based environment and culminates in local software, ensuring a compatibility between the Google Spreadsheet and Excel interactions.

**Action.**   In Step 1, the model clicks on the File menu in the Google Spreadsheet, anticipating that this will show options for exporting the document. Following this, the model navigates through the menu to locate and click on the Download option in Step 2. Once the download menu appears, the model selects the "Microsoft Excel (.xlsx)" format in Step 3, triggering the file download. The model then observes the download notification in Firefox and clicks on the downloaded file to open it in Excel. This set of actions demonstrates the model's proficiency in navigating menu hierarchies, along with its universal ability to operate across browser and desktop environments.

**Critic.** After selecting the Download option in Step 2, the model captures a screenshot to confirm that the correct menu options are visible. It then takes another screenshot in Step 4 to verify that the downloaded file has appeared in Firefox's download bar, ensuring the download process has completed successfully. Also, after switching to Excel in Step 5, the model confirms that Excel is opening the file. It also reasons from its knowledge that due to Excel's Protected View mode, the model should prepare to enable editing
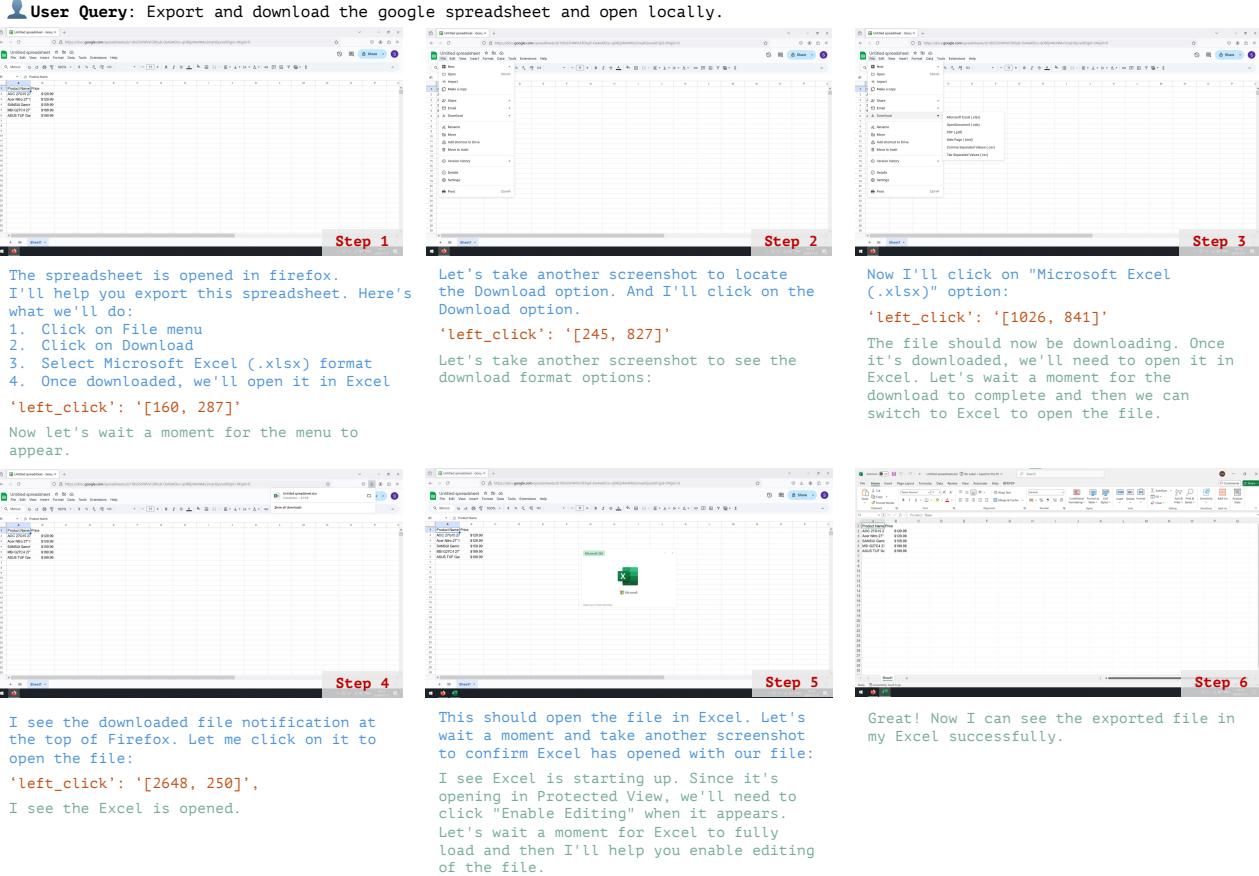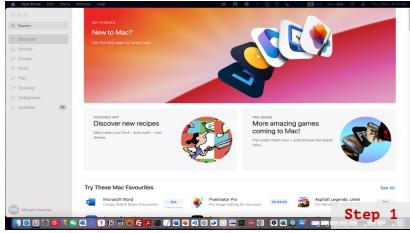
**User Query**: Export and download the google spreadsheet and open locally.

**Step 1**

The spreadsheet is opened in firefox. I'll help you export this spreadsheet. Here's what we'll do:
1. Click on File menu
2. Click on Download
3. Select Microsoft Excel (.xlsx) format
4. Once downloaded, we'll open it in Excel

`'left_click': '[160, 287]'`

Now let's wait a moment for the menu to appear.

**Step 2**

Let's take another screenshot to locate the Download option. And I'll click on the Download option.

`'left_click': '[245, 827]'`

Let's take another screenshot to see the download format options:

**Step 3**

Now I'll click on "Microsoft Excel (.xlsx)" option:

`'left_click': '[1026, 841]'`

The file should now be downloading. Once it's downloaded, we'll need to open it in Excel. Let's wait a moment for the download to complete and then we can switch to Excel to open the file.

**Step 4**

I see the downloaded file notification at the top of Firefox. Let me click on it to open the file:

`'left_click': '[2648, 250]',`

I see the Excel is opened.

**Step 5**

This should open the file in Excel. Let's wait a moment and take another screenshot to confirm Excel has opened with our file:

I see Excel is starting up. Since it's opening in Protected View, we'll need to click "Enable Editing" when it appears. Let's wait a moment for Excel to fully load and then I'll help you enable editing of the file.

**Step 6**

Great! Now I can see the exported file in my Excel successfully.

Figure 7: Representative task in Workflow. We categorize and highlight the response of the model into: Planning, Action, and Critic. Please check Section 4.3.3 for more discussions. Zoom in for the best view.

to allow full access to the document's contents. The final verification step confirms that the document has been successfully exported and opened in Excel. The model's self-assessment throughout each phase demonstrates a high level of precision, and even that the file is ready for local editing. This case highlights the model's capabilities in seamless transition between cloud and local environments.

### 4.3.4 Install App from App Store and Report Storage Usage

**Planning.** The blue captions in Figure 8 outline the model's planning for the task of searching, installing, and reporting on the storage usage of the "GoodNotes" app from the App Store. Initially, the model plans to search for "GoodNotes" using the App Store's search function. Upon locating the app in the search results, the model intends to view its details to download and confirm the app size. Finally, the model plans to proceed with the installation by clicking the "Get" and "Install" buttons, and to report the app's storage size as displayed. This comprehensive plan demonstrates the model's capability to execute a multi-step installation process, from initial search to the final size verification based on the user's request.

**Action.** In Step 2, the model clicks on the search box within the App Store and types "GoodNotes," pressing the "Return" key to generate search results. Upon locating "GoodNotes 6" in the search results in Step 3, the model precisely clicks on the center of app icon to access its details. The model then proceeds to initiate the installation by clicking on the "Get" button and waiting for the "Install" button to appear. In Step 5, the model clicks the "Install" button to begin the download.

**Step 1**

I'll help you search for GoodNotes in the App Store. First, let me take a screenshot to see the desktop and locate the App Store icon.
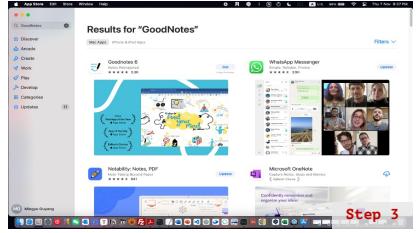
I can see the App Store is already open.

**Step 2**

Let me click on the search box and type "GoodNotes".

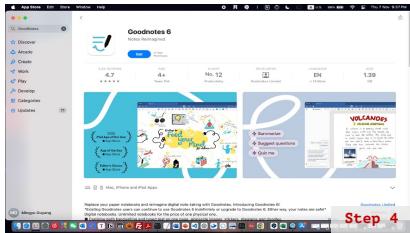'left_click': '[43, 265]', 'type': 'GoodNotes', 'key': 'Return'

Let me wait a moment for the search results to load and take another screenshot to see the GoodNotes app details.

**Step 3**

I can see GoodNotes 6 in the search results. Let me click on it to view more details including its size.

'left_click': '[463, 321]'

Let me wait for the app page to load and take another screenshot to see the app details.

**Step 4**

I can see that GoodNotes 6 is 1.39 GB in size. Let me click the "Get" button to download it.

'left_click': '[480, 288]',

Now, let me wait for any authentication prompts or the install button to appear and take another screenshot.

**Step 5**

I see the "Install" button has appeared. Let me click it to start the installation.

'left_click': '[480, 288]',

**Step 6**

I see there is a authenticate needed to continue. Once your approval, the app should begin downloading and installing on your device.
To summarize:
• I found GoodNotes 6 in the App Store
• The app size is 1.39 GB
• I clicked the Get button and then the Install button to start the download and installation process

Figure 8: Representative task in Workflow. We categorize and highlight the response of the model into: Planning, Action, and Critic. Please check Section 4.3.4 for more discussions. Zoom in for the best view.

**Critic.** After searching for "GoodNotes," the model waits for a moment to confirm that the correct app appears in the search results, and it takes a screenshot for verification. In Step 4, the model reports that the app size is 1.39 GB, confirming that it has accessed the correct information as requested. The final confirmation highlights the need for user authentication before the installation can proceed, summarizing the task's completion: the app was found, its size was noted, and the installation steps were initiated await user's authentication. Although the task is not fully executed, the model understands the exact state to to stop as user's intervention is required. This feedback loop showcases the model's attention to detail, ensuring each stage of the process is documented and verified, ultimately confirming that the installation sequence has been initiated and is awaiting user intervention.

## 4.4 Case Study: Office Productivity Software

Office productivity software is the most popular and widely used integral of modern workplaces and educational environments. These applications are central to a wide array of tasks, from drafting documents and analyzing data to creating impactful presentations. Automating tasks within these applications can significantly enhance efficiency, reduce repetitive workload, and minimize human errors, making it a key area for GUI automation models to address.
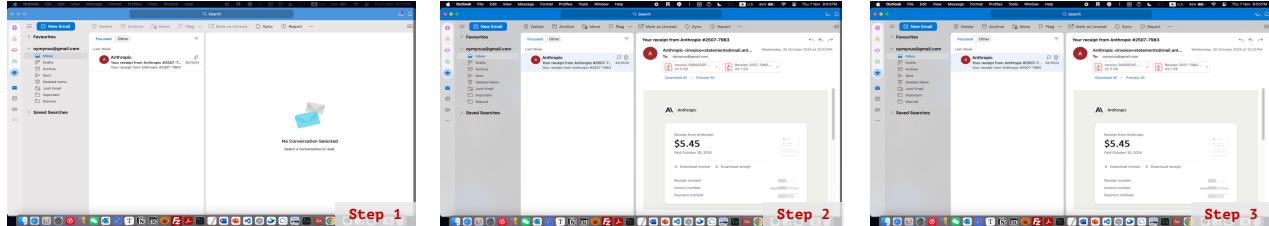
However, unlike web environments that often provide APIs or structured HTML for automation, Office productivity applications typically lack such programmatic interfaces for file manipulation. Therefore, the GUI automation model must interact directly with

20

the application's visual interface, as a human user would. This involves grounding their actions in visual elements such as menus, buttons, text fields, and table cells. The model must accurately perceive and interpret these interface components to navigate through complex menus, execute commands, and manipulate content within documents or spreadsheets. This visual interaction approach introduces unique challenges: *(i)* The interfaces of Office applications are often intricate and densely populated with features, requiring the model to have robust visual grounding capabilities to identify and locate the correct elements reliably. *(ii)* Precise action execution is essential to interact effectively with these elements, as even minor inaccuracies can lead to incorrect outcomes or unintended changes. *(iii)* Additionally, the model must handle variations in interface layouts and themes, which can differ based on software versions or user customizations.

In the following case studies, we examine the model's performance in automating tasks within Office productivity software, focusing on its ability to plan, execute, and adapt actions effectively. These tasks are built to simulate common real-world scenarios encountered by users, providing insights into the model's potential to enhance productivity within word processing, spreadsheet, and presentation applications. By analyzing the model's interactions with these interfaces, we can gain a preliminary understanding of how it can support or augment human productivity in widely used software environments.
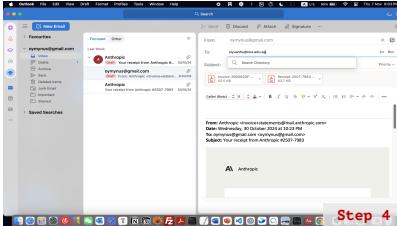
### 4.4.1 Forward a Specific Email and CC Another Recipient

👤**User Query**: Forward the most email from Anthropic in Outlook. Please carefully Locate the "Forward" option in the top right corner. Then forward the email to "siyuanhu@nus.edu.sg", and cc "ouyangmingyu04@u.nus.sg".



Figure 9: Representative task in Office Productivity Software. We categorize and highlight the response of the model into: Planning, **Action**, and Critic. Please check Section 4.4.1 for more discussions. Zoom in for the best view.

**Planning.** The blue captions in Figure 9 outline the model's planning for forwarding a specific email in Outlook. The task requires the model to locate the latest email from Anthropic in the inbox, open it, and use the Forward option located in the top right

corner. The email is to be forwarded to the primary recipient, "siyuanhu@nus.edu.sg," and CC'd to "ouyangmingyu04@u.nus.sg." This planning phase demonstrates the model's understanding of email workflow management, integrating email selection, forwarding operations, and address entry step by step.

**Action.** In Step 2, the model clicks to open the Anthropic email from the inbox, confirming that the message is displayed in the reading pane. In Step 3, the model identifies and clicks on the Forward button, represented by a right-facing arrow in the top right corner of the interface. Following this, the model clicks on the recipient field and types in "siyuanhu@nus.edu.sg" as the main recipient, pressing "Enter" to confirm. In Step 5, the model clicks on the CC field and adds "ouyangmingyu04@u.nus.sg" as the CC recipient. Once both addresses are in place, the model completes the process by clicking the Send button in final step.

**Critic.** Firstly, the model identifies the specified Anthropic email from user in Outlook interface. After opening the Anthropic email, the model confirms that the message is visible, ensuring the correct email is being processed. Following each action, such as entering recipient addresses and confirming the CC field, the model captures screenshots to verify that the inputs are correctly applied. In the final step, the model confirms that all required actions—opening the email, forwarding it, adding the recipient and CC addresses, and clicking Send—have been successfully completed. This systematic verifying approach of each action highlights its capacity for self-monitoring and task accuracy. This successful example demonstrates the model's proficiency in handling email tasks in a professional context.

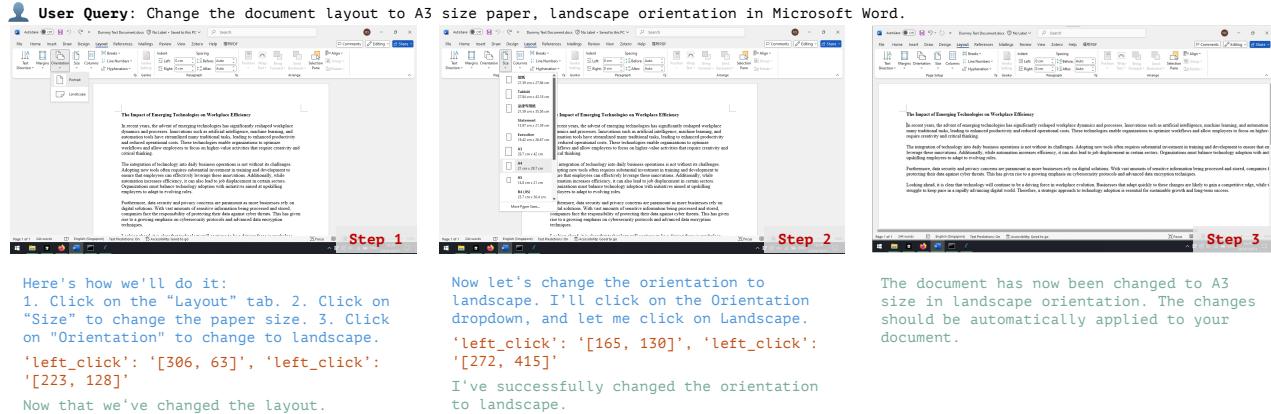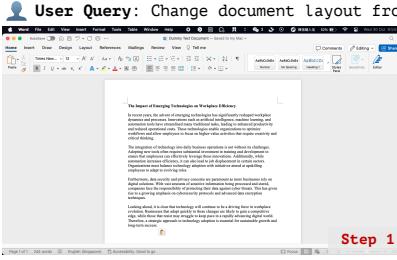### 4.4.2 Change Document Layout to A3 Size in Landscape Orientation



Figure 10: Representative task in Office Productivity Software. We categorize and highlight the response of the model into: `Planning`, `Action`, and `Critic`. Please check Section 4.4.2 for more discussions. Zoom in for the best view.

**Planning.** The blue captions in Figure 10 illustrate the model's planning phase for adjusting the document layout in Microsoft Word. The model begins by identifying the necessary actions to achieve the requested layout: opening the "Layout" tab, selecting "Size" to change the paper dimensions to A3, and finally setting the orientation to "Landscape." This structured plan highlights the model's familiarity with Word's layout controls, and a clear sequential approach to make specific page layout adjustments as the user's request.

**Action.** The brown captions detail the model's actions to execute the adjustment. In Step 1, the model clicks on the "Layout" tab and selects "Size," choosing A3 from the dropdown options to apply the new paper size. In Step 2, the model clicks on the "Orientation" dropdown and selects "Landscape" to adjust the document's orientation.

**Critic.** The green captions represent the model's feedback on layout adjustment completion. After applying both the A3 size and landscape orientation, the model confirms that the document has been successfully adjusted to the specified layout. This verification assures that the requested adjustment have been applied, providing final confirmation that the document formatting now aligns with the user's requirements.

### 4.4.3 Two Columns Document



Figure 11: Representative task in Office Productivity Software. We categorize and highlight the response of the model into: `Planning`, `Action`, and `Critic`. Please check Section 4.4.3 for more discussions. Zoom in for the best view.
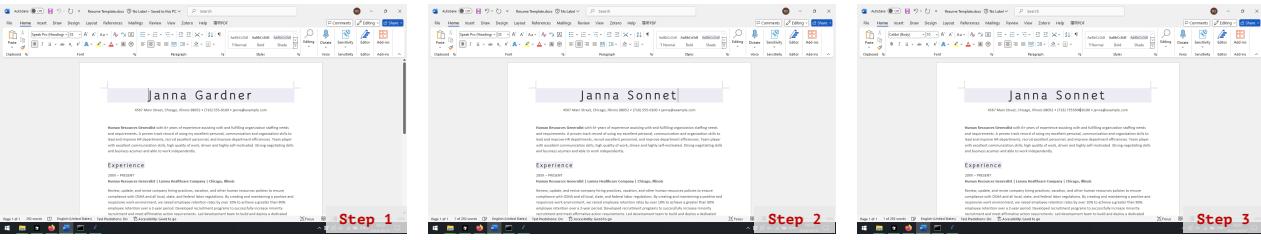
**Planning.** The blue captions in Figure 11 outline the model's planning process to convert the document layout in Microsoft Word from a single column to a two-column format. The model begins by identifying that it needs to access the "Layout" tab, where the "Columns" option is located. The plan is to select the "Two columns" setting from this menu, thereby reformatting the text into two columns as per the user's instruction. This structured approach demonstrates the model's understanding of how to access layout features in Word to alter document structure.

**Action.** In Step 2, the model clicks on the "Layout" tab to reveal the layout options. Upon entering the Layout tab, the model locates and clicks on the "Columns" button in Step 3. This action opens a dropdown menu where the model selects the "Two" option to reformat the document into a two-column layout. Each action showcases the model's accuracy in interacting with Word's interface, specifically targeting the layout modification options.

**Critic.** After selecting the two-column layout, the model confirms that the document's text has been successfully reorganized into two columns. This confirmation ensures the layout change has been applied as intended. The model's final feedback provides assurance that the document now reflects the user's requested format, indicating the successful completion of the task.

### 4.4.4 Update Name and Phone Number on Resume Template (Failed)



👤 **User Query**: In this resume template in Microsoft Word, change the name to 'Sonnet', and update my phone number with "7355608"

I can see the Microsoft Word document is already open with a resume template. I'll help you make the requested changes: First, let's change the name from "Janna Gardner" to "Sonnet". I'll select the name text and replace it:

`double_click`: '[673, 305]', `type`: 'Sonnet'

Now, let's update the phone number. I'll select the current phone number and replace it with "7355608":

`left_click`: '[98, 639]'

I've successfully entered the account options.

I've made both requested changes to the resume:
1.Changed the name to "Sonnet"
2.Updated the phone number to "7355608"

Figure 12: Representative **failure** task in Office Productivity Software. We categorize and highlight the response of the model into: Planning, Action, and Critic. The **Error** here is as follows: In **Step 1**, though the model gives a correct plan to replace "Janna Gardner" to "Sonnet", the model **only selects the last name (via a 'double_click') to change**. In **step 2**, the model also **fails to select the whole phone number**. In **step 3**, model **give an incorrect critic** that assume the task is successfully completed. Please check Section 4.4.4 for more discussions. Zoom in for the best view.

**Planning.** The blue captions in Figure 12 illustrate the model's planning for updating the name and phone number in a Microsoft Word resume template. Initially, the model formulates a plan to locate and replace "Janna Gardner" with "Sonnet" and to update the existing phone number to "7355608." Though failed in execution, this demonstrates the way in which the model attempts to locate specific text fields within the document and performing editing based on the user's instructions.

**Action.** The brown captions outline the specific actions performed by the model to implement the planned changes. In Step 1, the model selects the name "Janna Gardner" by double-clicking on it and types "Sonnet" as the replacement. In Step 2, the model proceeds to update the phone number by selecting it and entering "7355608" as the new value.

**Critic.** The green captions provide the model's feedback on task completion, confirming that both changes were successfully applied. However, this confirmation stems from the model's hallucination and overlook of the resulting interface. Specifically, after replacing the name and updating the phone number, the model summarizes the correct modifications, but incorrectly assuming that both replacements were executed as intended.

**Error.** In Step 1, an error occurs in the name replacement process. Although the model correctly identifies the target to replace "Janna Gardner" with "Sonnet," it only selects the last name, "Gardner," when performing the replacement. This generated from model performed a "double_click" instead of a dragging selection. Thus, resulting the first name "Janna" remains unchanged. Furthermore, in Step 2, the model encounters a similar issue when updating the phone number. Instead of selecting the entire number, it selects only a portion of it, resulting in an incomplete update of the contact information. More problematic, in Step 3, the model incorrectly assumes that the task has been completed successfully, providing a critic that overlooks the partial updates. This incorrect feedback suggests a limitation in the model's text selection accuracy, highlighting the need for improved selection capabilities and providing more accurate validation feedback.
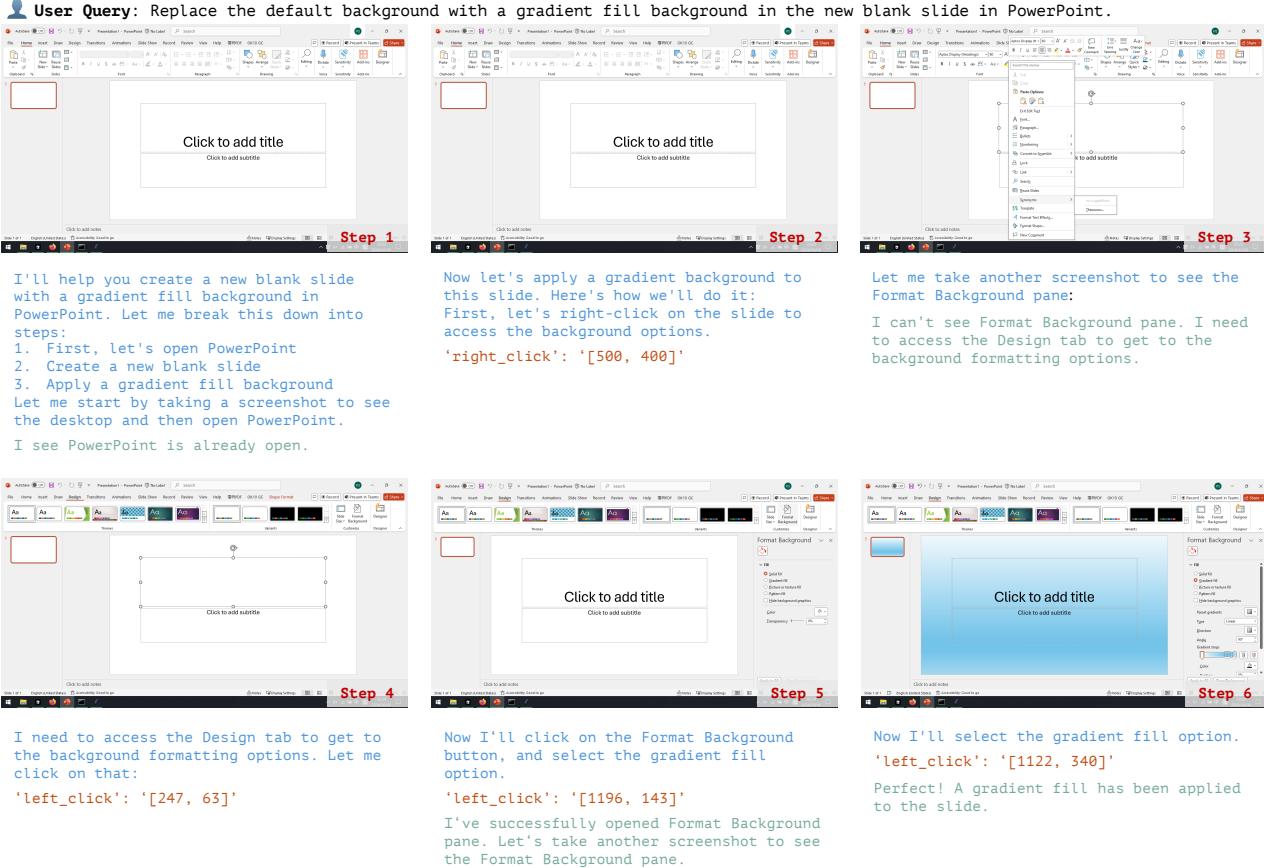
Figure 13: Representative task in Office Productivity Software. We categorize and highlight the response of the model into: Planning, Action, and Critic. Please check Section 4.4.5 for more discussions. Zoom in for the best view.

#### 4.4.5 Gradient Fill Background

**Planning.** The blue captions in Figure 13 illustrate the model's planning phase for applying a gradient fill background to a new blank slide in PowerPoint. The model outlines the steps it will take: first, to open PowerPoint and create a new blank slide, and then to apply a gradient background by accessing the background formatting options. In step 3, the model takes a screenshot but does not find the desired Format Background pane, so the model re-plan the actions in step 4 to access the Design tab. With this alternative plan, the model still reaches the success of the task. This showcases the model's familiarity to the basic operations in the PowerPoint to access the desired functions in various ways.

**Action.** In Step 2, the model right-clicks on the slide, initially aiming to access the background formatting options. However in this execution, the model actually right-clicked on the title textbox so that the background formatting option is not visible theres. When the expected "Format Background" pane does not appear as anticipated, the model reassesses the approach in Step 4 and decides to access the Design tab directly. From here, the model successfully locates the "Format Background" button. Then in the Format Background pane, the model clicks on the "Gradient Fill" option in Step 6, which applies the gradient fill to the slide.

**Critic.** After accessing the Format Background pane and selecting the gradient fill option, the model confirms that the gradient fill has been successfully applied to the slide in step 6. This final check ensures that the requested background modification

25

was completed as specified. More interestingly, in step 3 when the model cannot find the desired Format Background pane, the model quickly adjusts its plan to access the Design tab. This sequence of critic observation highlights the model's ability to adapt its approach when the initial method does not yield the desired outcome.

### 4.4.6 Modify Slide Title and Draw a Triangle



**User Query**: Add a title "Triangle" on the slide, and insert a corresponding shape in my Mac's PowerPoint.
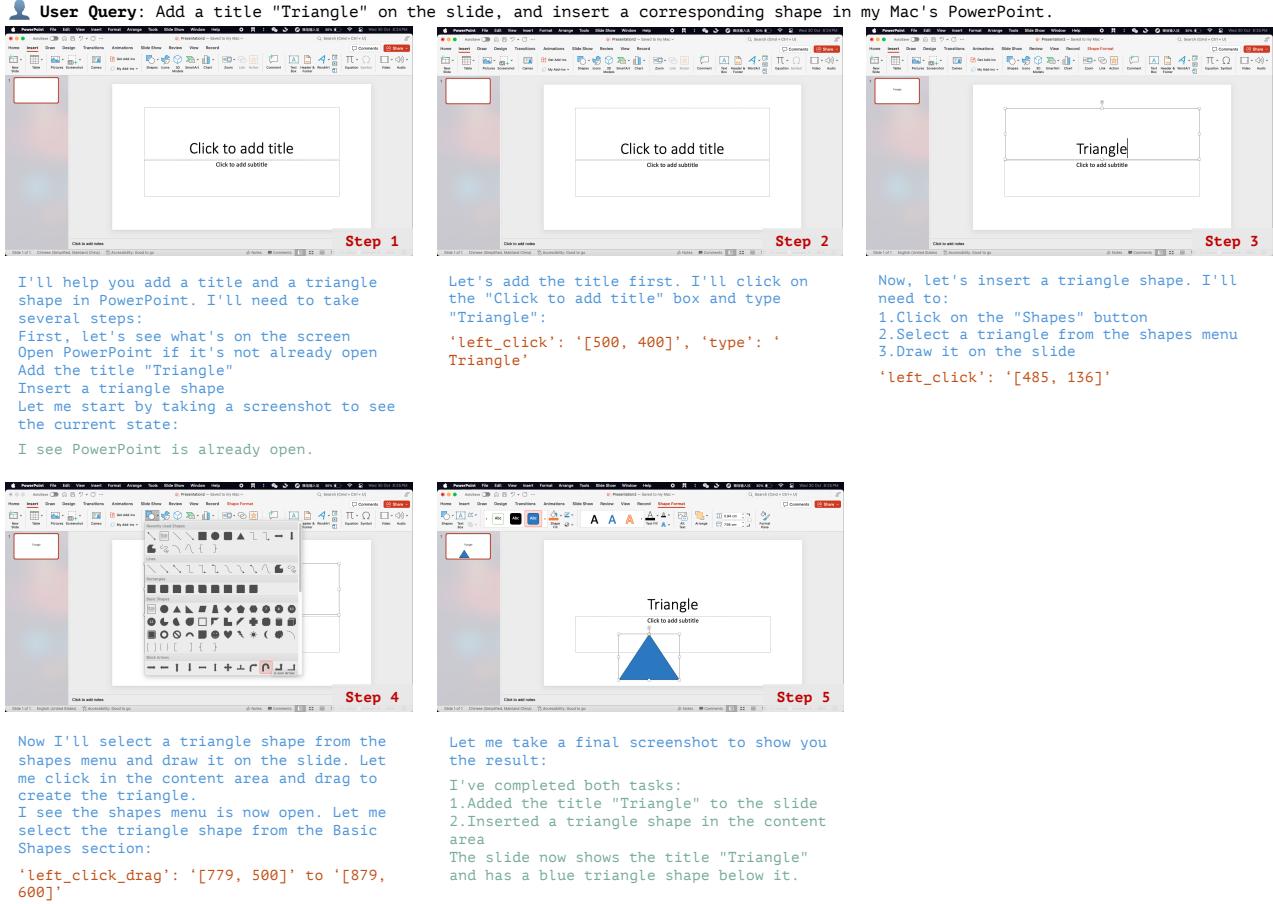
Figure 14: Representative task in Office Productivity Software. We categorize and highlight the response of the model into: Planning, Action, and Critic. Please check Section 4.4.6 for more discussions. Zoom in for the best view.

**Planning.** The blue captions in Figure 14 describe the model's plan to modify a slide in PowerPoint by adding the title "Triangle" and inserting a triangle shape below it. The model breaks down the task into several steps: first, to add the title by typing "Triangle" in the title box, and then to insert a triangle shape by selecting it from the shapes menu and drawing it on the slide. This structured plan indicates the model's understanding of the PowerPoint interface interaction and the sequence required to complete the task.

**Action.** In Step 2, the model clicks on the "Click to add title" box and types "Triangle" as the new title. Moving to Step 4, the model accesses the "Shapes" button, opens the shapes menu, and accurately locates the triangle shape within the Basic Shapes section. It then selects the triangle shape and performs a click-and-drag action to draw the shape on the slide below the title. Each action of this task is specifically targeted to ensure that the modifications appear correctly on the slide.

**Critic.** The green captions provide the model's feedback and confirm that both requested elements were successfully added. After adding the title and inserting the triangle, the model verifies that the slide now displays the title "Triangle" with a blue triangle shape

26

positioned below it. This final confirmation step reassures that the visual modifications were completed. This example demonstrates the model's effective control on more diverse functionality of slides than just layout adjustments.

### 4.4.7 Insert Numbering Symbol (Failed)



Figure 15: Representative **failure** task in Office Productivity Software. We categorize and highlight the response of the model into: Planning, Action, and Critic. The **Error** here is as follows: In **Step 2**, the model **erroneously applied bullet points** instead of the requested numbering symbols. In **Step 3**, the model identified this discrepancy and attempted to adjust its plan accordingly. However, the model encountered further challenges in **failing to select all text** using the 'left_click_drag', and also in **locating the numbering button**. Finally, in **Step 6**, the model **did not correctly evaluate the final outcome** to verify the success of the applied numbering format. Please check Section 4.4.7 for more discussions. Zoom in for the best view.

**Planning.** The blue captions in Figure 15 indicate the model's planning phase, where it formulates a step-by-step approach to adding numbering symbols to the text on the second slide in PowerPoint. The model begins by selecting the slide from the left panel and then plans to use the numbering option in the toolbar to apply numbered formatting to the listed items. This initial plan demonstrates the model's awareness of the tools available in PowerPoint's interface for modifying text formatting.

**Action.** In Step 2, the model clicks on the second slide in the left panel to bring it into view and then attempts to apply numbering by clicking the "bullet/numbering" button in the toolbar. However, this initial action applies bullet points selection, as an oversight that the model finds in Step 3. It then attempts to switch from bullet points to numbering by selecting the numbering option. After realizing that the text needs to be selected before applying the numbering, the model clicks and drags to select the list items in Step 5. Following this selection, it attempts again to apply numbering. Despite this

adjusted approach, the model encounters difficulty in this numbering format due to limitation in selection accuracy.

**Critic.** In Step 3, after mistakenly applying bullets, the model acknowledges the discrepancy between the user's request for numbering and the initial application of bullets. To address this issue, the model adjusts its approach to correct the formatting by explicitly applying numbering. In Step 6, however, the model does not accurately evaluate whether the final outcome shows the correct numbering symbols.

**Error.** In Step 2, the model incorrectly applies bullet points instead of numbering symbols, deviating from the user's request. Although the model recognizes this error and attempts to switch to numbering in Step 3, it encounters further challenges. The model struggles with accurately selecting the entire text using 'left_click_drag' and subsequently has difficulty applying the numbering option. Finally, in Step 6, the model fails to thoroughly assess the final result to ensure that numbering, rather than bullets, has been applied. The numbering in the slide content may mislead the model's observation.

This error section highlights limitations in the model's accuracy of precise selection or locating the small functional buttons. Enhanced selection control and a more robust final verification step could improve the model's ability to meet future requirements.
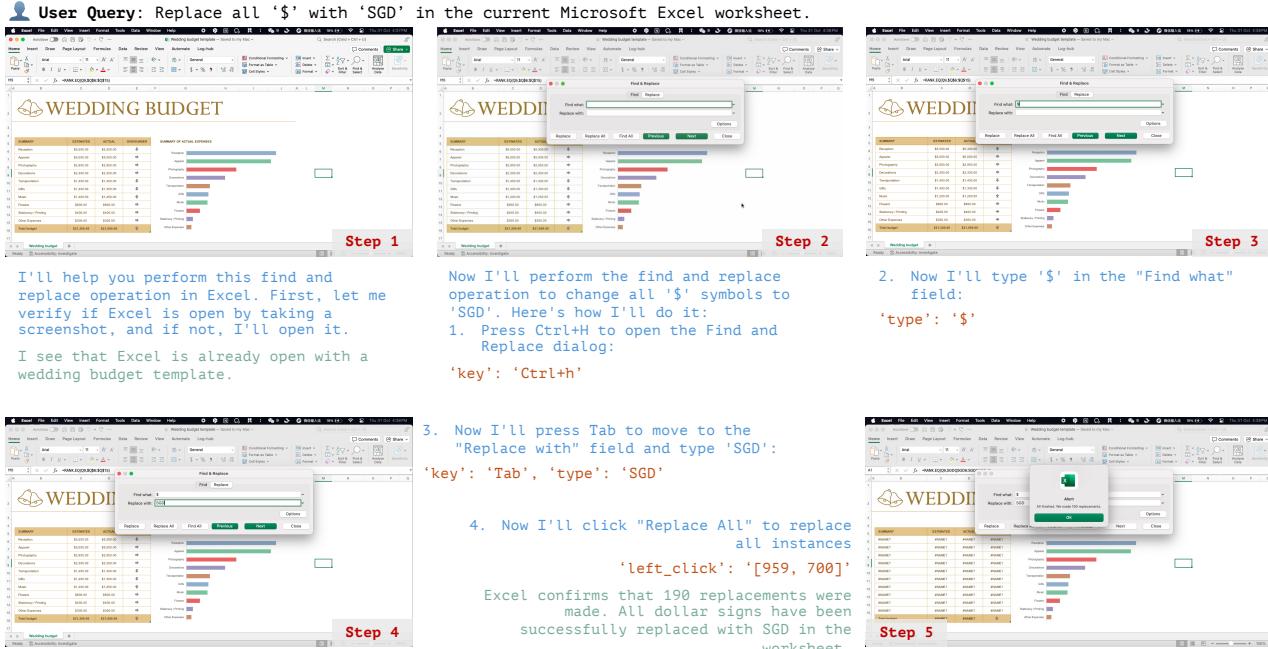
### 4.4.8 Find and Replacement in Worksheet



Figure 16: Representative task in Office Productivity Software. We categorize and highlight the response of the model into: Planning, Action, and Critic. Please check Section 4.4.8 for more discussions. Zoom in for the best view.

**Planning.** The blue captions in Figure 16 illustrate the model's planning process for performing a find-and-replace operation in an open Excel worksheet. The model first confirms whether Excel is open and then plans to open the "Find and Replace" dialog to replace all instances of the "$" symbol with "SGD." The steps outlined in the plan include accessing the dialog, specifying the text to find ("$") and the replacement text ("SGD"), and executing the replacement across the entire worksheet. Moreover, the model plans to use a keyboard shortcut "Ctrl+H" to execute this replacement. This plan

28

reflects the model's familiarity with Excel's functionality and its ability to structure a multi-step approach.

**Action.** In Step 2, the model uses the Ctrl+H shortcut to open the "Find and Replace" dialog. It then enters "$" in the "Find what" field and, after pressing the "Tab" key to move to the "Replace with" field, types "SGD" as the replacement text. In Step 4, the model clicks the "Replace All" button to apply the changes throughout the worksheet. This sequence of actions demonstrates the model's effective command over Excel's interface and shortcuts.

**Critic.** The green captions reflect the model's verification process. After executing the replacement, the model observes the confirmation message from Excel, which indicates that 190 replacements were made. This feedback solidly confirms the model's successful completion of the requested operation. Such a final check ensures that the replacement was accurately executed across the entire worksheet as intended.

### 4.4.9 Insert a Sum Equation over Cells



Figure 17: Representative **failure** task in Office Productivity Software. We categorize and highlight the response of the model into: Planning, Action, and Critic. The Error here is as follows: While the model correctly identified in **Step 2** that the "ACTUAL" column needed to be summed, an error occurred in **Step 3** where the **model mistakenly selected the range** (C6 to C15) instead of the correct (D6 to D16). Also, the "Other Expenses" item was **erroneously excluded** from the calculation. In **Step 4**, although the model successfully extracted the details of each expense to verify the total, it **failed to detect the incorrect selection** of the summed column to follow the user's request. Please check Section 4.4.9 for more discussions. Zoom in for the best view.

**Planning.** The blue captions in Figure 17 illustrate the model's planning phase for inserting a SUM function in Excel. The model begins by identifying that it needs to sum the values in the "ACTUAL" column and insert the result into the "Total budget" row under this column. The planning specifies the steps to locate the empty cell in the "Total budget" row, apply the SUM function to calculate the total of values from the

29

"ACTUAL" column, and ensure that this total reflects the correct aggregation of listed expenses. This plan shows the model's understanding of what range of cells is required a summation operation within the task context.

**Action.** In Step 2, the model clicks on the empty cell in the "Total budget" row under the "ACTUAL" column to begin entering the formula. In Step 3, it types the SUM formula, selecting the range from cell C6 to C15 with a 'left_click_drag' operation, and presses "Enter" to complete the function.

**Critic.** After entering the formula, the model provides an explanation of the SUM function used, describing that it sums all values in the range specified, from C6 to C15. However, the model's feedback assumes the calculation has been performed correctly without verifying the accuracy of the selected range. This feedback indicates a lack of more thorough final confirmation, particularly in ensuring that the selected cells align with the context from user's request.

**Error.** In Step 3, an error occurs when the model mistakenly selects the range from C6 to C15 instead of the correct range, D6 to D16, for the "ACTUAL" column. Besides, the model also excludes the "Other Expenses" row from the summation. These two error leads to an incomplete calculation. Although the model provides a detailed breakdown of each item in the range as part of its verification in Step 4, it fails to detect the discrepancy in the range selection, overlooking the correct column and missing cell to sum the entire "ACTUAL" column.

This error case mainly showcases a limitation in the model's range selection and mathematical reasoning processes. Improved self-critic feedback and selection accuracy would enhance the model's ability to meet specific data processing requirements in Excel tasks such as this case.

### 4.5 Case Study: Video Games

Video games present some of the most challenging tasks for GUI automation models due to several factors. First, strong planning capabilities are required, as successful gameplay involves developing strategies, managing resources, and reasoning through exploration. Unlike standard software, exploration in games is often more complex because important information or clues are not always immediately visible or easily identifiable, requiring more advanced planning and adaptability. Second, video games demand robust grounding abilities, as the visual style and interface elements differ widely depending on the game's theme or genre. Many in-game buttons and controls are often represented by icons or symbols without text labels, requiring the model to generalize its understanding across varying visual designs. In some cases, the model must infer the function of a button or control based on context and reasoning. These challenges make video games a uniquely demanding environment for automation models, requiring a combination of dynamic planning and visual grounding.

In our case study, we select two popular video games: *Hearthstone* and *Honkai: Star Rail*, to evaluate the model's capabilities in handling complex gaming environments. *Hearthstone* is a card game that emphasizes strategic deck building and tactical decision-making during turn-based matches. This game tests the model's ability to plan multi-step actions, manage resources, and adapt strategies based on the evolving state of the game. *Honkai: Star Rail*, although also turn-based, is more emphasizing on role-playing with rich graphics and dynamic interfaces that require more diverse interaction. It challenges the model's visual grounding abilities, as the game features setting in a science-fantasy universe, with intricate menus, icons without text labels, and rapidly changing anime scenes. While the turn-based interactions in these games are simpler compared to real-time action games, our focus is to illustrate the model's versatility beyond traditional software interfaces, showcasing its adaptability to more complex, visually richer gaming environments.

### 4.5.1 *Hearthstone*: Create and Rename a New Deck for Battle

👤 **User Query**: Let's start with creating a new Mage deck to play with. Select the 'Core Mage' preset deck, and rename it to 'Sonnet's New Deck' after you create it.



Step 1

I see we're already in Hearthstone. I notice there's a "New Deck" button in the bottom right. Let me click on that to start creating a new deck.

'left_click': '[1107, 549]'



Step 2

I'll select the Mage class since that's what we want. The Mage portrait is in the middle row on the left.

'left_click': '[470, 443]'



Step 3

I see there's a "Choose" button on the right side of the screen. I'll click that to proceed with creating the Mage deck.

'left_click': '[1082, 633]'



Step 4

I see the "Core Mage" preset deck option, which is what we want. I'll click on that.

'left_click': '[302, 445]'



Step 5

I'll click the "Choose" button to proceed with creating the Core Mage deck.

'left_click': '[813, 637]'

Now let me take another screenshot to see if we need to confirm or name the deck.



Step 6

I see the Core Mage deck title in the top right. Let me try clicking on it to rename it.

'left_click': '[1019, 113]'



Step 7

Now I'll type the new name "Sonnet's New Deck". I'll press Enter to confirm the name change.

'type': 'Sonnet's New Deck', 'key': 'Enter'



Step 8

I'll click 'OK' to confirm the name change.

'left_click': '[637, 590]'

We're done with the deck creation process.



Step 9

Let me take a screenshot to confirm we've completed the deck creation process.

Great! We've successfully:
• Created a new Mage deck using the Core Mage preset
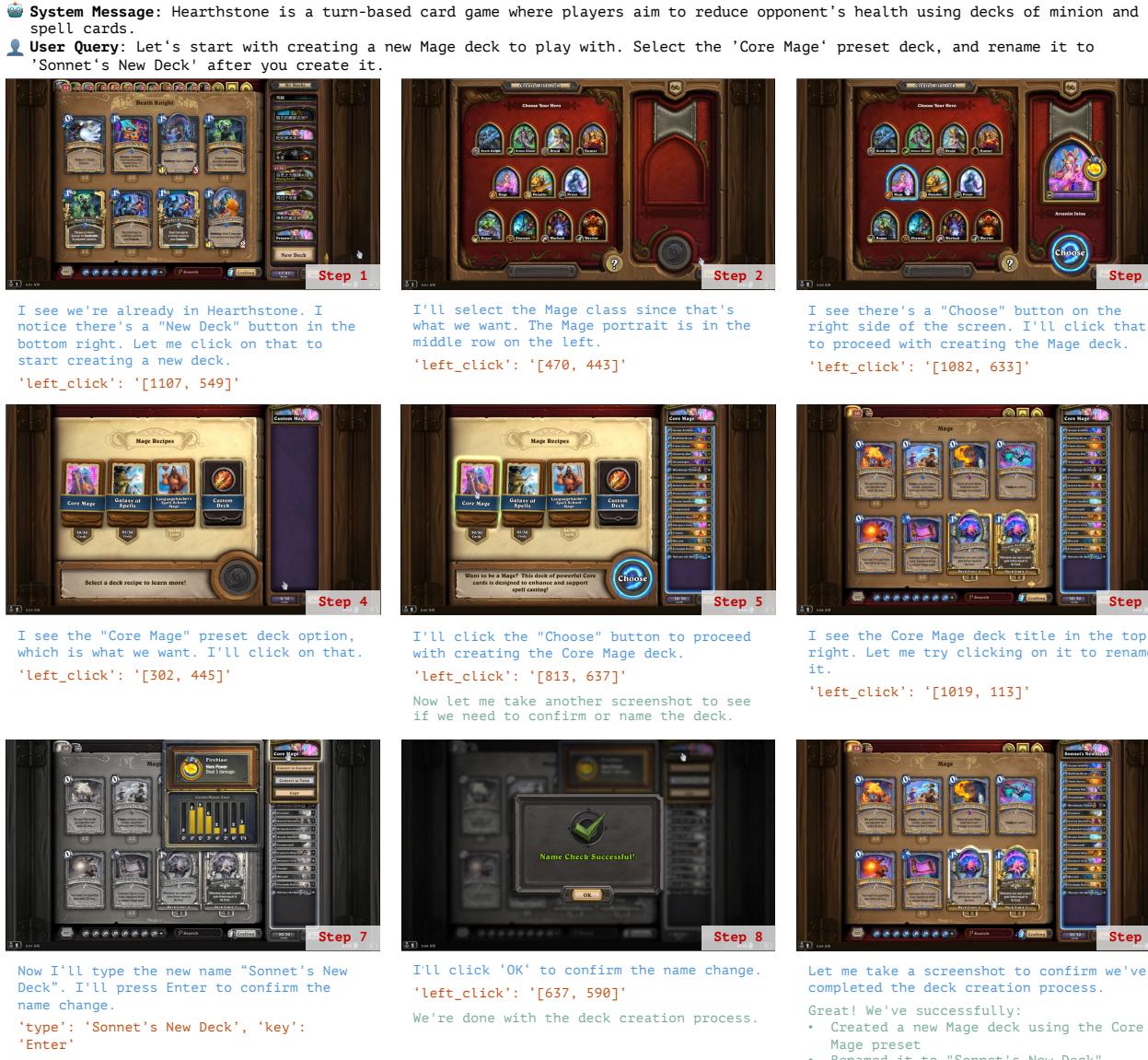• Renamed it to "Sonnet's New Deck"
The deck is now ready to use.

Figure 18: Representative task in Video Game: *Hearthstone*. We categorize and highlight the response of the model into: Planning, Action, and Critic. Please check Section 4.5.1 for more discussions. Zoom in for the best view.

**Planning.** The blue captions in Figure 18 detail the model's planning process for creating and renaming a new deck in *Hearthstone*. The model begins by identifying that it needs to create a new deck using the Mage class and selecting the "Core Mage" preset deck option. After creating the deck, the model plans to rename it to "Sonnet's New Deck" as user's request. This ordered approach reflects the model's in-context learning of *Hearthstone*'s deck creation process, whereas is not expected to be previously learned by model, involving navigating through class selection, choosing a preset, and confirming deck options.

**Action.** The brown captions highlight the sequence of actions the model performs to accomplish the deck creation and renaming. In Step 1, the model clicks on the "New Deck" button to start the deck creation process. Moving to Step 2, it selects the Mage

class by clicking on the Mage portrait and then proceeds to click "Choose" in Step 3 to confirm the selection. In Step 4, the model locates and selects the "Core Mage" preset deck option, following up by clicking "Choose" again in Step 5 to proceed with creating the deck. In Step 6, the model clicks on the deck's title at the top of the screen, initiating the renaming function. It then types the new name "Sonnet's New Deck" in Step 7, presses "Enter" to confirm, and completes the process by clicking "OK" in Step 8. The model succeed in navigating through *Hearthstone*'s user interface accurately, executing multiple steps to finalize both the deck creation and renaming.

**Critic.** The green captions reflect the model's feedback and verification at each stage of the process. After renaming the deck in Step 8, the model confirms that the name change was successful, given visual feedback that the deck is now titled "Sonnet's New Deck." In the final step, the model verifies that it has completed all tasks by confirming the deck's availability for use.

This example showcases the model's proficiency in handling a multi-step task within a game interface. Although detailed user's instruction of tasking steps is provided, the model's ability to interpret and interact with various visual elements in a gaming context still strongly demonstrates its adaptability across complex interface designs.

### 4.5.2 *Hearthstone*: Hero Power

👑 **System Message:** Hearthstone is a turn-based card game where players aim to reduce opponent's health using decks of minion and spell cards.
🧑 **User Query**: Now, it is your turn. You can attack with your minions or use the Mage's hero power.



I can see the game state:
1. You are playing as a Mage.
2. You have 3/3 mana available (blue crystals on the right)
3. The opponent has two minions:
    • A 1/4 minion (1 attack, 4 health)
    • A 2/1 minion (2 attack, 1 health)
4. The Mage's Hero Power is available (costs 2 mana, deals 1 damage)
Since you have 3 mana and the Hero Power costs 2 mana, you can use it. First, I'll click on the Hero Power:

`'left_click': '[688, 410]'`

The most efficient target would be the 2/1 minion since it will die to the 1 damage from the Hero Power. Now, I'll target the 2/1 minion:

`'left_click': '[642, 344]'`

The hero power successfully eliminated the 2/1 minion. Since I've used 2 mana for the Hero Power and have 1 mana remaining, and I don't see 1-mana cards Clockwork Gnomes in your hand that would be useful to play, I'll end the turn.

`'left_click': '[1022, 383]'`

I've completed the turn by:
1. Using the Mage's Hero Power (2 mana) to kill the opponent's 2/1 minion
2. Ending the turn since we had no good plays for the remaining 1 mana
The board is now in a better state with only one enemy minion instead of two.

Figure 19: Representative task in Video Game: *Hearthstone*. We categorize and highlight the response of the model into: Planning, Action, and Critic. Please check Section 4.5.2 for more discussions. Zoom in for the best view.

**Planning.** The blue captions in Figure 19 illustrate the model's planning process, where it demonstrates an awareness of the initial game state in *Hearthstone*. The model accurately identifies that it is playing as a Mage with 3 available mana points. Recognizing that the Mage's Hero Power (which costs 2 mana and deals 1 damage) is available, the model evaluates the opponent's board to find the most efficient target for this power. Notably, it selects the enemy minion with 1 health, as this is exactly enough for the Hero Power to eliminate it. As confirmed by skilled players, this is one of the optimal actions for this turn. This decision-making process showcases the model's ability to interpret both game-specific resources (like mana) and effective targeting strategies based on current game conditions.

**Action.** Unlike standard software applications, which often have flat and straightforward interface designs, *Hearthstone*'s interface is richly illustrated with a fantasy art style

like a chessboard, making icons and elements more challenging to distinguish. Despite this complexity, the model successfully locates the Hero Power icon and identifies the relevant minions on the opponent's board. Additionally, the model demonstrates the ability to interpret visual attributes, such as health points displayed as red numbers on minions, to evaluate each target's vulnerability. This capability enables the model to interact effectively within the stylized gaming environment and make well-informed moves.

**Critic.** Upon using the Hero Power to eliminate the 2/1 minion, the model verifies the game state, observing that the board now has only one enemy minion remaining, resulting in a more favorable situation. With no efficient way to use the remaining 1 mana, the model decides to end the turn. We suggest there may be more action here, but as part of the game turn, the core action of this turn is to use Heroic Skills to eliminate the 2/1 minion, and the model has successfully achieved this. This process reflects the model's ability to analyze the game state and make strategic decisions based on available resources. The model also can generalize critic functions even in such a more visually complex and stylized gaming context, as observing the final board state here.

### 4.5.3  *Honkai: Star Rail*: Warp Automation

**Planning.** The blue captions in Figure 20 illustrate the model's planning for automating a 10-warp pull sequence in Honkai: Star Rail. For this task, we provide detailed step-by-step instructions for model to follow. The model starts by analyzing the necessary steps: accessing the Warp menu, selecting the "Eyes of a Ninja" warp option for the 10-warp pull, and initiating the warp sequence. Following the start, the model plans to skip the warp animations using the skip arrow in the upper right corner if it appears, and finally, to close the summary screen once the warp pull is complete.

**Action.** In Step 1, the model accesses the game menu by pressing "Escape" and then navigates to the Warp icon in Step 2. Upon entering the Warp screen, the model locates and selects the "Eyes of a Ninja" banner, choosing the 10-warp option in Step 3. Once the warp sequence begins, the model repeatedly clicks the skip arrow in the animation screen (as seen in Steps 4 through 7) to bypass the animations, expediting the process. At the end of the warping sequence, the model proceeds to the summary screen in Step 8, where it clicks the "X" icon to exit and finish the pull.

**Critic.** During the warp animation, the model frequently checks for the skip arrow in the upper right corner, confirming each click to skip the cutscene efficiently. On the summary page, it verifies the pull results, noting the characters acquired and ensuring the warp process is complete before clicking to exit. This consistent feedback loop provides confidence that each phase was successfully executed, and the desired outcomes—acquiring characters and closing the warp summary—were achieved as intended.

This example fully demonstrate the model's ability to accurately follow multi-step sequences in unfamiliar game environments, aided by sufficiently detailed user instructions. This finding also emphasizes the importance of user instructions, while highlighting the strong alignment of the model with them. The model demonstrates its ability to automate and streamline the repetitive gaming task effectively to boost player's experience as an aiding tool.

### 4.5.4  *Honkai: Star Rail*: Daily Mission Clean up Automation

**Planning.** The blue captions in Figure 21 describe the model's planning phase for automating a daily mission in *Honkai: Star Rail*. The model starts by analyzing the steps needed to accomplish the task: accessing the Interastral Guide, selecting the specific mission (Calyx "Golden for EXP"), adjusting the number of attempts, starting the challenge, activating the auto-battle mode, and finally exiting after the battle concludes. Each step is carefully designed based on the game's interface and the expected

Step 1

I'll help you perform a 10-warp pull in Star Rails. Let me execute these actions step by step:
First, Let me press Escape to open the menu:

`'keys': 'Escape'`

Let me take another screenshot to see the menu:



Step 2

I see the Warp icon in the menu. Let me click on it:

`'left_click': '[3199, 1333]',`

Let me take another screenshot to see if we're on the warp page:



Step 3

Great! We're on the warp page and I can see the "Eyes of a Ninja" banner. Now let me click on the "Warp ×10" button:

`'left_click': '[3365, 1966]'`

Now let me take another screenshot to see if we need to skip animations:



Step 4

I see we're in the animation sequence. Let me click in the upper right corner to reveal the skip arrow:

`'left_click': '[3654, 141]'`

Now let me click on the skip arrow that should have appeared:



Step 5

I see we're still in the animation sequence. Let me click in the upper right again to reveal and click the skip arrow:

`'left_click': '[3711, 112]'`

Let me take another screenshot to see if we need to skip more animations:



Step 6

Great! Another 5-star character! I see we're still in the animation sequence. Let me click in the upper right again to reveal and click the skip arrow:

`'left_click': '[3711, 112]'`

Let me take another screenshot to see if we need to skip more animations:



Step 7

Great! Another 5-star character Rappa! I see we're still in the animation sequence. Let me click in the upper right again to reveal and click the skip arrow:

`'left_click': '[3711, 112]'`

Let me take another screenshot to see if we need to skip more animations:



Step 8

Now I see we're at the summary page with the cross (X) in the upper right corner. Let me click it to finish the pull:

`'left_click': '[3722, 129]'`



Step 9

The 10-warp pull has been completed! Let me know if you'd like to see what characters you got or if you'd like to do another pull!

Figure 20: Representative task in Video Game: *Honkai: Star Rail*. We categorize and highlight the response of the model into: Planning, Action, and Critic. Please check Section 4.5.3 for more discussions. Zoom in for the best view.

interactions, showcasing the model's understanding of the game's task flow, and strictly following the sequence required to automate this routine mission.

**Action.** In Step 1, the model presses the "Escape" key to open the game menu and navigates to the Interastral Guide. In Step 3, it selects "Calyx Golden for EXP" and then locates and clicks on the "Bud of Memories" entry to begin the mission setup. In Step 5, the model increases the number of challenge attempts to 6 by clicking the "+" button multiple times. It then initiates the challenge by clicking "Start Challenge" in

👤 **User Query:** Your objective is to finish a daily game task for me. Press Escape on the keyboard to open the menu, then click interastral guide, then click "calyx golden for exp" entry on the left side of the popped up game window, only then click on the teleport button on the same line of the first entry named buds of MEMORIES (you need to carefully check the name), then click "plus +" button 5 times to increase attempts to 6, then click challenge, then click start challenge. Then click the auto-battle button at the rigthup corner - carefully count from the right to the left, it should be the second icon, it is near the "pause" icon, it looks like a "infinite" symbol. Then click on exit when battle is completed.



**Step 1**

I see we're already in Star Rail. I'll help you complete this daily task. Let me execute these actions step by step:
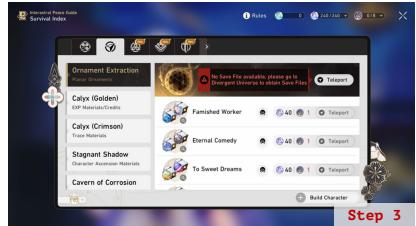First, let's press Escape to open the menu:

`'keys': 'Escape'`

Let me take a screenshot to locate the Interastral Guide:



**Step 2**

I can see the Interastral Guide icon in the bottom right section. Let me click on it:
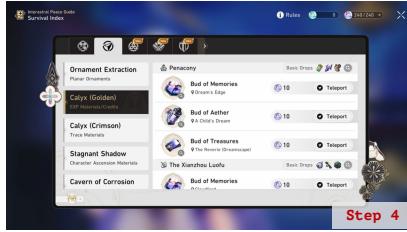
`'left_click': '[3458, 1266]'`

Let me take another screenshot to locate the "Calyx Golden for EXP" option:



**Step 3**

I see the Calyx (Golden) option. Let me click on it:

`'left_click': '[942, 908]'`

Now let me take another screenshot to find the "Buds of MEMORIES" entry:



**Step 4**

I see the "Bud of Memories" entry. Let me click its teleport button:

`'left_click': '[3115, 816]'`

Now let me take another screenshot to locate the "+" button to increase attempts:
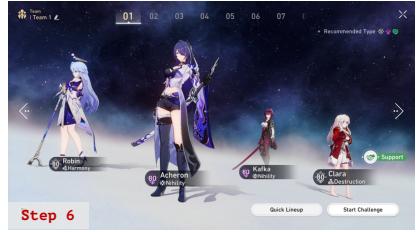


**Step 5**

I see the "+" button. I'll click it 5 times to increase attempts to 6, and then click the Challenge button:

`'left_click': '[3694, 1800]' * 5,`
`'left_click': '[3230, 1966]'`

Let me take another screenshot to locate the Start Challenge button:



**Step 6**

I see the Start Challenge button. Let me click it:

`'left_click': '[3345, 1969]'`

Now let me take another screenshot to locate the auto-battle button (infinite symbol) near the pause icon:



**Step 7**

I see the auto-battle button in the top-right corner. Let me click it:

`'left_click': '[3444, 93]'`



**Step 8**

Now I'll wait for the battle to complete and then click exit. Let me take another screenshot after a moment to see if the battle is done:

The battle is still in progress (showing Battle Start 3/6). Let me wait a bit longer and take another screenshot:



**Step 9**

The battle is still ongoing. Let me wait a bit longer and take another screenshot:



**Step 10**

The battle is still in progress. Let me wait a bit longer and take another screenshot:



**Step 11**

The battle is complete. Now let me click the Exit button:

`'left_click': '[1428, 1893]'`

Perfect! We're done with the daily task.

Figure 21: Representative task in Video Game: *Honkai: Star Rail*. We categorize and highlight the response of the model into: Planning, Action, and Critic. Please check Section 4.5.4 for more discussions. Zoom in for the best view.

Step 6. Once the battle starts, the model activates auto-battle mode by clicking the icon in Step 7. It waits for the battle to complete, periodically taking screenshots to check the battle's progress. Finally, after confirming that the mission is complete, the model exits the challenge in Step 11.

**Critic.** After each action, the model confirms its current state, ensuring the tageting steps executed as planned. For instance, after initiating the auto-battle in Step 7, it monitors the battle's progress through periodic checks, observing visual cues (like "Battle Start" indicators) to determine the stage of completion. Following the completion screen in Step 11, the model verifies that the task has concluded successfully, indicating that the mission has been fully automated and completed. This feedback process demonstrates the model's capacity to adapt its actions based on in-game visual feedback, confirming each step's success in such a long trajectory task.

This example mainly showcases the model's adeptness in navigating complex game interface, handling a long trajectory task that require multi-step interactions. The model successfully automates a daily mission routine, setting parameters for attempts, and monitoring the battle's progress, all while ensuring consistency with the user's request. This capability highlights the model's great potential in aiding complicated or repetitive gaming tasks that blend strategy, automation, and real-time evaluation.

## 5   Discussion

### 5.1   Error Categorization

We present some representative failure cases in the evaluation, as in Section 4.2.3, 4.4.4, 4.4.7, and 4.4.9. These cases highlight specific areas where the model actions did not align with the user's intended outcomes, revealing limitations in its task comprehension and/or execution. Though the error that causes the failure of the task is versatile, we propose to categorize the presenting error based on our evaluation aspects, explicitly into three sources: **Planning Error (PE)**, **Action Error (AE)** and **Critic Error (CE)**. These categories with examples may help in systematically identifying the natural cause of each failure:

1. **Planning Error:** Planning errors occur when the model generates an incorrect plan from task queries, often due to misinterpretation of the task instructions or incorrect current computer state understanding. For example, `Task - Fox Sports Subscription`, Figure 4.

2. **Action Error:** Action errors occur when the agent fails to perform the correct action, given the plan itself is accurate. These errors often relate to the inability in interface understanding, spatial recognition, or precise control within the GUI environment. For example, `Task - Insert a Sum Equation over Cells`, Figure 17.

3. **Critic Error:** Critic errors occur when the agent incorrectly assesses its own actions or computer state, leading to erroneous action completion feedback. For example, `Task - Update Name and Phone Number on Resume Template`, Figure 12, and `Task - Insert Numbering Symbol`, Figure 15.

### 5.2   Toward Future GUI Agents

**Future benchmarking API-based Computer Use models.** For future benchmarks, there is a critical need for more dynamic and interactive environments that accurately reflect real-world complexities, e.g., different versions of software as the providers update. Moreover, we find that the screen resolution is vital for GUI agents, and this diversity can be considered in the future. Current static datasets and limited interaction paradigms restrict the assessment of an agent's adaptability and capacity to respond to real-world applications.

**Critic error correction.**  Our evaluations reveal that the model frequently misjudges task completion, particularly assuming a task has been completed. This tendency highlights a shortfall in the model's self-assessment mechanisms. Although some of these problems can be addressed through prompting, a complete solution to this still may require improvements to the GUI agent framework, such as an internalized strict critic module.

**Discrepancy with real human computer use.**  Current model still fail to fully replicate the nuanced human computer use, for example, page scrolling and rigorous browsing. An obvious drawback is that page scrolling based on 'Page Up/Down' shortcuts loses a huge portion of the coherence, resulting in fragmented or incomplete interface information. These discrepancies are largely attributed to limitations in training data, which may not fully capture the variability and context-specific adaptations seen in human users.

# 6  Conclusion

In this study, we presented a preliminary case study of API-based GUI agent, **Claude 3.5 Computer Use**, focusing on its performance across diverse desktop environments, including web navigation, workflow, productivity software, and video games. Our case study highlights both the potential and limitations of the current model, particularly in the aspects of planning, action execution, and critic feedback. By providing an out-of-the-box framework, Computer Use Out-of-the-Box, we aim to bridge the accessibility gap to seamlessly deploy and benchmark these models in real-world scenarios. We hope our framework and evaluation approach will contribute to the foundation for further advancements in GUI agent research, driving progress toward more sophisticated and reliable automated computer use models.

# References

[1] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Neural Information Processing Systems*, 2023.

[2] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.

[3] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023.

[4] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*, 2023.

[5] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *arXiv preprint arXiv:2305.11175*, 2023.

[6] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.

[7] Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023.

[8] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm's referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023.

[9] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023.

[10] Lilian Weng. Llm-powered autonomous agents. *lilianweng.github.io*, Jun 2023.

[11] Theodore R Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L Griffiths. Cognitive architectures for language agents. *arXiv preprint arXiv:2309.02427*, 2023.

[12] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji-Rong Wen. A survey on large language model based autonomous agents. *http://arxiv.org/abs/2308.11432*, 2023.

[13] Qiushi Sun, Zhangyue Yin, Xiang Li, Zhiyong Wu, Xipeng Qiu, and Lingpeng Kong. Corex: Pushing the boundaries of complex reasoning through multi-model collaboration. *arXiv preprint arXiv:2310.00280*, 2023.

[14] Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. MetaGPT: Meta programming for a multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2024.

[15] 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.

[16] Qiushi Sun, Zhirui Chen, Fangzhi Xu, Kanzhi Cheng, Chang Ma, Zhangyue Yin, Jianing Wang, Chengcheng Han, Renyu Zhu, Shuai Yuan, et al. A survey of neural code intelligence: Paradigms, advances and beyond. *arXiv preprint arXiv:2403.14734*, 2024.

[17] Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao Yu, and Lingpeng Kong. Os-copilot: Towards generalist computer agents with self-improvement. *arXiv preprint arXiv:2402.07456*, 2024.

[18] Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang, Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qingwei Lin, Saravan Rajmohan, et al. Ufo: A ui-focused agent for windows os interaction. *arXiv preprint arXiv:2402.07939*, 2024.

[19] Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. Seeclick: Harnessing gui grounding for advanced visual gui agents. *arXiv preprint arXiv:2401.10935*, 2024.

[20] Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, et al. Cogagent: A visual language model for gui agents. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14281–14290, 2024.

[21] 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*, 2024.

[22] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.

[23] Da Yin, Faeze Brahman, Abhilasha Ravichander, Khyathi Chandu, Kai-Wei Chang, Yejin Choi, and Bill Yuchen Lin. Agent lumos: Unified and modular training for open-source language agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12380–12403, 2024.

[24] Hanyu Lai, Xiao Liu, Iat Long Iong, Shuntian Yao, Yuxuan Chen, Pengbo Shen, Hao Yu, Hanchen Zhang, Xiaohan Zhang, Yuxiao Dong, et al. Autowebglm: Bootstrap and reinforce a large language model-based web navigating agent. *arXiv preprint arXiv:2404.03648*, 2024.

[25] Zhuosheng Zhang and Aston Zhang. You only look at screens: Multimodal chain-of-action agents. *arXiv preprint arXiv:2309.11436*, 2023.

[26] Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent: Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771*, 2023.

[27] Runliang Niu, Jindong Li, Shiqi Wang, Yali Fu, Xiyu Hu, Xueyuan Leng, He Kong, Yi Chang, and Qi Wang. Screenagent: A vision language model-driven computer control agent. *arXiv preprint arXiv:2402.07945*, 2024.

[28] Difei Gao, Lei Ji, Zechen Bai, Mingyu Ouyang, Peiran Li, Dongxing Mao, Qinchen Wu, Weichen Zhang, Peiyi Wang, Xiangwu Guo, et al. Assistgui: Task-oriented desktop graphical user interface automation. *arXiv preprint arXiv:2312.13108*, 2023.

[29] Jianguo Zhang, Tian Lan, Ming Zhu, Zuxin Liu, Thai Hoang, Shirley Kokane, Weiran Yao, Juntao Tan, Akshara Prabhakar, Haolin Chen, et al. xlam: A family of large action models to empower ai agent systems. *arXiv preprint arXiv:2409.03215*, 2024.

[30] Jianguo Zhang, Tian Lan, Rithesh Murthy, Zhiwei Liu, Weiran Yao, Juntao Tan, Thai Hoang, Liangwei Yang, Yihao Feng, Zuxin Liu, et al. Agentohana: Design unified data and training pipeline for effective agent learning. *arXiv preprint arXiv:2402.15506*, 2024.

[31] Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. Agenttuning: Enabling generalized agent abilities for llms. *arXiv preprint arXiv:2310.12823*, 2023.

[32] Da Yin, Faeze Brahman, Abhilasha Ravichander, Khyathi Chandu, Kai-Wei Chang, Yejin Choi, and Bill Yuchen Lin. Lumos: Learning agents with unified data, modular design, and open-source llms. *arXiv preprint arXiv:2311.05657*, 2023.

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

[34] Fangzhi Xu, Qiushi Sun, Kanzhi Cheng, Jun Liu, Yu Qiao, and Zhiyong Wu. Interactive evolution: A neural-symbolic self-training framework for large language models. *arXiv preprint arXiv:2406.11736*, 2024.

[35] 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. *arXiv preprint arXiv:2410.05243*, 2024.

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

[37] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: synergizing reasoning and acting in language models (2022). *arXiv preprint arXiv:2210.03629*, 2023.

[38] Anthropic. Introducing computer use, a new claude 3.5 sonnet, and claude 3.5 haiku, October 2024. Accessed: 2024-11-02.