

PC Agent: While You Sleep, AI Works - A Cognitive Journey into Digital World

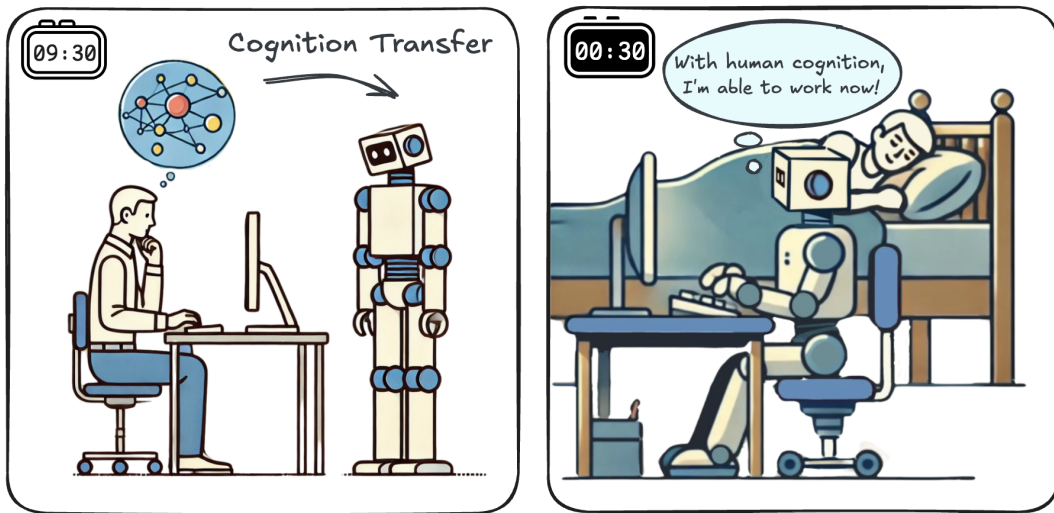
Yanheng He^{1,3*} Jiahe Jin^{1,3*}
Shijie Xia^{1,2,3} Jiadi Su³ Runze Fan^{1,3} Haoyang Zou³ Xiangkun Hu³ Pengfei Liu^{1,2,3†}

¹Shanghai Jiao Tong University ²SII
³Generative AI Research Lab (GAIR)

Abstract

Imagine a world where AI can handle your work while you sleep - organizing your research materials, drafting a report, or creating a presentation you need for tomorrow. However, while current digital agents can perform simple tasks, they are far from capable of handling the complex real-world work that humans routinely perform. We present **PC Agent**, an AI system that demonstrates a crucial step toward this vision through **human cognition transfer**. Our key insight is that the path from executing simple “tasks” to handling complex “work” lies in efficiently capturing and learning from human cognitive processes during computer use. To validate this hypothesis, we introduce three key innovations: (1) PC Tracker, a lightweight infrastructure that efficiently collects high-quality human-computer interaction trajectories with complete cognitive context; (2) a two-stage cognition completion pipeline that transforms raw interaction data into rich cognitive trajectories by completing action semantics and thought processes; and (3) a multi-agent system combining a planning agent for decision-making with a grounding agent for robust visual grounding. Our preliminary experiments in PowerPoint presentation creation reveal that complex digital work capabilities can be achieved with a small amount of high-quality cognitive data - PC Agent, trained on just 133 cognitive trajectories, can handle sophisticated work scenarios involving up to 50 steps across multiple applications. This demonstrates the data efficiency of our approach, highlighting that the key to training capable digital agents lies in collecting human cognitive data. By open-sourcing our complete framework, including the data collection infrastructure and cognition completion methods, we aim to lower the barriers for the research community to develop truly capable digital agents. Resources are available at <https://gair-nlp.github.io/PC-Agent/>.

arXiv:2412.17589v1 [cs.AI] 23 Dec 2024



Human work, AI learn

AI work, Human rest!

*Equal Contribution.

†Corresponding author.

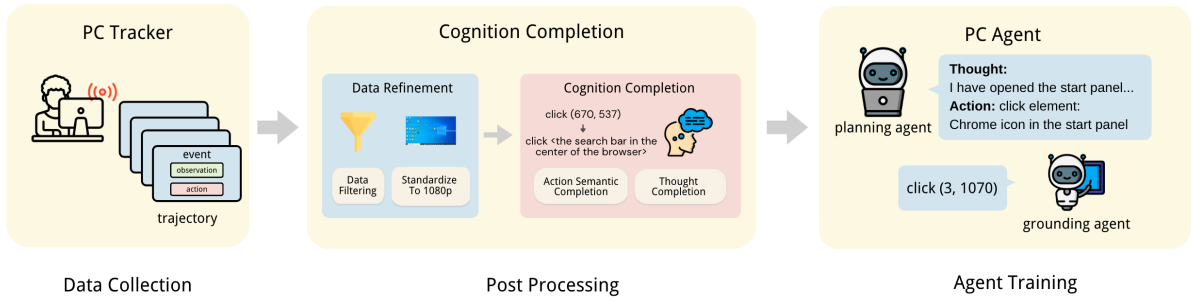


Figure 1: Overview of our framework, consisting of three key components: (1) PC Tracker, a lightweight infrastructure that collects human-computer interaction trajectories by recording user actions and state observations; (2) a two-stage cognition completion that converts raw interaction data into cognitive trajectories through data refinement and human cognition completion, including action semantics and thought processes; and (3) a multi-agent system comprising a planning agent for action decision-making and a grounding agent for click position grounding.

1 Introduction

While artificial intelligence (AI) has made remarkable progress in understanding and generating content like text (OpenAI, 2022, 2024a) and images (Rombach et al., 2022; Ramesh et al., 2022; Li et al., 2022), these capabilities exist largely in the representation world, where AI systems master pattern recognition and symbolic processing (Bengio et al., 2014). Recent advances in large language model (LLM) powered autonomous agents have garnered significant attention for their ability to perform task-oriented interactions (Weng, 2023; Wu et al., 2023). In the domain of digital agents - AI systems that assist with digital tasks - current solutions have demonstrated basic capabilities in simple tasks like web searches (Yao et al., 2023a; Deng et al., 2023). However, they remain far from capable of handling complex real-world work that humans routinely perform (Xie et al., 2024; Bonatti et al., 2024), such as video editing, presentation creation, or report writing - activities that require sustained operation across multiple applications, sophisticated decision-making, and even aesthetic judgment. To meaningfully reduce human workload, digital agents must evolve beyond simple task execution to manage more complex work efficiently.

The Critical Gap: From Simple Task Automation to Complex Work Completion Most existing approaches in building digital agents rely heavily on proprietary LLM APIs (Wu et al., 2024a; Zheng et al., 2024; Zhang et al., 2023). While few works have attempted to train models for computer operation, these efforts have primarily concentrated on either improving basic visual grounding capabilities (Gou et al., 2024; Wu et al., 2024b) or specific domains like web (Liu et al., 2024). However, these approaches still struggle with complex real-world computer work. Through our analysis, we identify two critical challenges: (1) **foundational visual grounding**: The ability to precisely locate GUI elements (e.g. the Start button in the taskbar) represents a fundamental capability that current vision-language models still struggle with (Xie et al., 2024; Gou et al., 2024). (2) **complex cognitive understanding**: More crucially, current agents lack the cognitive capabilities required for complex work. They struggle with maintaining contextual awareness across extended interaction sequences, making dynamic decisions based on changing environments, and adapting strategies in response to execution outcomes. While prompt engineering (Liu et al., 2021) can partially help, it proves inadequate for truly complex work.

Human Cognition: The Missing Key Our key insight is that the path to the digital world lies in **human cognition transfer**. When completing complex work, the human brain engages in sophisticated cognitive activities - understanding objectives, analyzing current states, reflecting on past actions, and planning future strategies. These cognitive processes ultimately crystallize into clear decisions, which are then externalized into observable behavior. This observation leads us to develop a novel framework that efficiently captures and transfers human cognition to AI agents.

Our Contributions In this work, we present three key contributions that together form a complete solution for developing truly capable digital agents:

- **PC Tracker**: We introduce the first open-source infrastructure for efficiently collecting large-scale human-computer interaction trajectories. This lightweight system captures not just actions but the complete cognitive context of human computer use, establishing a rich foundation for agent training.
- **Cognition Completion Pipeline**: We develop an innovative approach that transforms raw interaction data into cognitive trajectories by completing action semantics and thought processes in steps. This enables AI to learn not just what humans do, but why and how they make decisions in complex digital environments.

-
- **PC Agent:** We propose an AI system that can perform complex cognitive work in the digital world. Using PowerPoint presentation creation as our main testing ground, PC Agent shows remarkable capabilities in executing long-sequence tasks and switching between applications, allowing it to create complex presentations with practical utility.

Findings and Impact Our preliminary experiments demonstrate **significant data efficiency** in learning from human cognitive processes rather than just behavioral data - PC Agent, trained on just 133 cognitive trajectories, can execute complex work with up to 50 steps. By **open-sourcing our framework**, we provide the research community with more than just a proof of concept; we offer a practical foundation for further research and development. As we progress toward a broader vision - from representation world to digital world, and ultimately to the physical world - this work represents a critical step in enabling AI to practically enhance human productivity.

2 How Far Are We From True Digital Agents

2.1 The Capability Gap in Real-world Applications: From “task” to “work”

In the evolution of digital agents, LLM-powered autonomous agents have garnered significant attention due to their ability to engage in complex task-oriented interactions (Weng, 2023; Wu et al., 2023). Early works like ReAct (Yao et al., 2023b) and Reflexion (Shinn et al., 2023) established fundamental frameworks for agent reasoning and self-optimization. In implementing digital agents, two primary technical approaches emerged: backend access (Trivedi et al., 2024; Team, 2024) and frontend GUI interaction. The latter has gained prominence by operating without backend access permissions, offering universal applicability, better security, and adaptability to interface changes. With the advancement of vision-language models (Liu et al., 2023; OpenAI, 2024b,c; Wang et al., 2024), GUI agents have made significant progress in both visual understanding (Hong et al., 2023; Liu et al., 2024; You et al., 2024; Yang et al., 2023; Cheng et al., 2024; Lu et al., 2024) and general task capabilities, demonstrated by applications across web browsers (Gur et al., 2024; Zheng et al., 2024; Yao et al., 2023a; Nakano et al., 2022), mobile devices (Zhang et al., 2023; Hoscilowicz et al., 2024; Zhang et al., 2024b; Zhang and Zhang, 2024) and desktops (Wu et al., 2024a; Zhang et al., 2024a). To advance the field, researchers have established important evaluation benchmarks through comprehensive datasets (Rawles et al., 2023; Deng et al., 2023) and execution-based environments that simulate real-world scenarios (Zhou et al., 2024; Koh et al., 2024; Xie et al., 2024; Bonatti et al., 2024; Rawles et al., 2024).

Despite notable progress in the field, we remain far from achieving truly capable digital agents. Even the most advanced systems, including the recently celebrated new Claude-3.5-Sonnet (Anthropic, 2024), still significantly under-perform humans in computer use (Xie et al., 2024; Hu et al., 2024). While current digital agents can perform simple **tasks** like web searches and file copying, they face significant challenges when tackling comprehensive **work** that better reflect real-world computer use. Consider video editing, presentation creation, and report generation - these work require sustained operation across multiple applications, sophisticated decision-making, and even human-level aesthetic judgment. We argue that true digital agents should be able to efficiently handle these complex work, not just simple tasks, to meaningfully reduce human workload.

2.2 Breaking Barriers: Visual Grounding and Cognitive Understanding

To realize true digital agents, we have identified two critical technical challenges: foundational visual grounding capabilities and deep cognitive understanding.

Visual Grounding Visual grounding - the ability to precisely locate elements in GUI - represents a fundamental capability for agents to effectively use computers. This is because click-based interactions, which require precise coordinate outputs, constitute a significant portion of computer tasks. However, most current vision-language models (VLMs), including state-of-the-art proprietary models like GPT-4o and open-source models like Qwen2-VL, still lack this basic capability. Contemporary works such as UGround (Gou et al., 2024) and OS-ATLAS (Wu et al., 2024b) attempt to address this issue by fine-tuning models on large-scale GUI visual grounding datasets. While successful in developing the target capability, these approaches compromise the model’s general question-answering and instruction-following abilities. Notably, Anthropic’s new Claude-3.5-Sonnet has emerged as the first frontier model with state-of-the-art computer use capabilities, though its training details remain undisclosed. We believe that both specialization and general capabilities are equally important for effective agent models. Encouragingly, we discover that the open-source general-purpose VLM Molmo (Deitke et al., 2024) demonstrates exceptional visual grounding abilities, establishing a strong foundation for open-source digital agents.

Cognitive Understanding However, even with robust visual grounding capabilities, current agents still struggle to complete complex computer work. Our analysis reveals that cognitive understanding represents the crucial missing piece in this puzzle - one that we will address through our human cognition transfer framework. This limitation manifests in two key aspects:

- **Lack of Fine-grained Cognitive Knowledge:** While current agents can perform high-level task planning, they lack the fine-grained cognitive knowledge for specific GUI operations. For instance, adding a title in PowerPoint requires clicking the title box before typing. Such operational sequences, intuitive to humans, present challenges for agents. Moreover, when handling cross-application tasks, such as collecting images from browsers to presentations, agents must comprehend the interaction logic and division of labor between different applications. This challenge necessitates reconsidering agent training from a human cognitive perspective, as GUIs were fundamentally designed for human interaction.
- **Insufficient Agent Training:** While current language models exhibit strong capabilities in following instructions and generating responses, they often struggle with the unique demands of agent workflows. These workflows require maintaining contextual awareness across extended interaction sequences, making dynamic decisions based on changing environments, and adapting strategies in response to execution outcomes. Although prompt engineering can partially bridge this gap, it proves inadequate for truly complex work. This misalignment between model capabilities and agent requirements highlights the need for specialized training approaches to enable effective real-world work completion.

3 Cognition Transfer: A Gateway for AI into the Digital World

The journey to the digital world represents a fundamental evolution from understanding and responding in language to taking meaningful action for LLMs. Although current digital agents can perform simple tasks, they are still far from effectively completing real-world work in digital environments, where humans manage their everyday work and life activities. The fundamental challenge lies here: *how can we enable AI systems to fundamentally evolve from understanding tasks to executing them in the digital world?*

We posit that **human cognition transfer** is pivotal in addressing this challenge. Human cognition is primarily based on the dynamic interaction between internal thinking and external behavior. When performing complex work in real environments, the human brain engages in sophisticated cognitive activities, including understanding objectives, analyzing current states, reflecting on the past, and planning future strategies. This series of cognitive processes ultimately crystallizes into clear decisions, which are then externalized into observable behavior. In this sense, behavior serves as the external projection of complex cognitive activities in the human brain.

If AI systems can acquire human cognition to interact with the digital world, they can complete complex work as naturally as humans do. However, current technology cannot directly record cognitive activities in the brain. To address this, we propose an indirect approach:

1. First, we designed a lightweight infrastructure to **efficiently collect raw interaction trajectories** between humans and digital devices. These trajectories not only capture the specific steps of task execution but also reflect dynamic exploration, trial and error, and optimization.
2. Then, we leverage LLMs to analyze these raw trajectories and **complete the cognitive processes** behind human behavior by completing action semantics and thought processes in steps, transforming them into *cognitive trajectories*. These cognitive trajectories can be considered effective approximations of genuine human cognitive activities.

By learning from these cognitive trajectories, AI systems can not only mimic specific operational behaviors but also master the underlying human cognition more efficiently. Our experimental results demonstrate remarkable data efficiency in enabling AI systems to perform complex computer work through human cognition transfer.

4 PC Tracker: Human-Computer Interaction Data Collection Infrastructure

While supervised fine-tuning has proven effective for adapting LLMs to many tasks, large-scale computer interaction trajectories remain severely limited compared to conventional text or image datasets. This data bottleneck has become a key impediment in developing digital agents that can effectively use computers like humans do.

To address this bottleneck, we present PC Tracker, the first lightweight infrastructure for efficient large-scale collection of human-computer interaction data. Similar to screen recording mechanisms, PC Tracker runs seamlessly in the background, recording user actions by monitoring keyboard and mouse activities while capturing screenshots to document state observations. This enables the collection of real-world interaction trajectories (see an example in Figure 3) at scale, establishing a rich foundation for future research in agent training and cognition engineering. PC Tracker is designed with the following key features, as summarized in Figure 2.

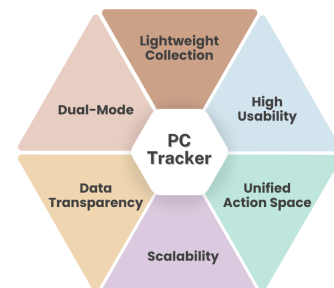


Figure 2: Key features of PC Tracker

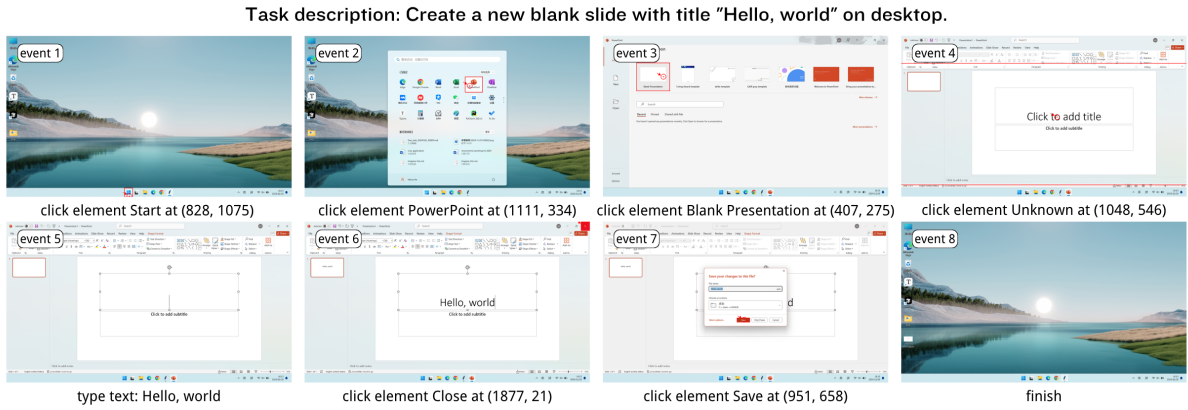


Figure 3: An example trajectory collected by PC Tracker. Red marks on the screenshots indicate the positions of click-related actions.

Lightweight Collection Unlike bulky and redundant video recording solutions, PC Tracker efficiently collects data through event-based tracking. It captures critical events $e = \langle act, obs \rangle$ when user operation is detected, where act represents action and obs represents the corresponding computer state observation. Critical events compose the trajectory, capturing the complete and authentic human-computer interaction while significantly reducing storage usage.

High Usability Running seamlessly in the background, PC Tracker allows users to use their computers naturally. Unlike some contemporary approaches that record complete accessibility trees, we deliberately avoid this practice as the crawling process would introduce noticeable latency and disrupt user operations.

Scalability The lightweight design and high usability ensure the feasibility of large-scale, long-term deployment, enabling PC Tracker to support unlimited-scale data collection. Our statistics show that an hour of continuous computer usage typically generates around 2,000 events, demonstrating immense data potential.

Unified Action Space We encapsulate fragmented keyboard and mouse operations into a unified action space. Raw operations are combined into actions like double click and type, significantly reducing the number of recorded events while enhancing action semantics.

Data Transparency We ensure user privacy through local data storage, complete recording control, and transparent Markdown visualization of all stored trajectories.

Dual-Mode PC Tracker offers two recording modes: task oriented and non-task oriented. The non-task oriented mode captures user interaction trajectories without specific tasks, ideal for large-scale data collection, while the task oriented mode primarily serves for supervised fine-tuning data annotation.

4.1 Tracking Strategy

4.1.1 Action Recording

The original records of certain actions are often fragmented, leading to a loss of semantic information, such as double-click, scroll, hotkey, and type. For instance, a type action is split into individual key press operations. Besides, some actions require special rules to be identified from the original operations, such as the *drag to* action, which is originally just a normal mouse release.

To address this issue, we designed heuristic algorithms that track the history of mouse and keyboard activities to encapsulate these raw operations into a unified action space \mathcal{A} designed based on OSWorld (Xie et al., 2024), as shown in Figure 4. This action space, originally designed for AI agents, significantly reduces the difficulty for AI to understand and learn human operational behaviors. For example, Figure 5 demonstrates how PC Tracker encapsulates 9 raw actions into a single *type text* action when the user intends to type "Hello". The user's raw actions involve capitalization change, character typing, and error corrections, as shown in the left column. These raw actions are

Action	Description
<i>click (x, y)</i>	clicks at coordinates.
<i>right click (x, y)</i>	right-click at coordinates.
<i>double click (x, y)</i>	double-click at coordinates.
<i>press (x, y)</i>	press mouse down at coordinates.
<i>drag to (x, y)</i>	drags the mouse to coordinates.
<i>scroll (0, 10)</i>	scrolls the screen with offset dy = 10.
<i>press key: enter</i>	presses the Enter key.
<i>hotkey (ctrl, c)</i>	performs the Ctrl+C hotkey (copy).
<i>type text: hello</i>	type text "hello".
<i>wait</i>	pauses for some time.
<i>finish</i>	the task is finished.
<i>fail</i>	the task is failed.

Figure 4: Action space \mathcal{A} of PC Tracker.

accumulated into the type buffer and combined into the unified action “type text: Hello” upon detecting type completion.

Additionally, we observed that coordinate-based click-related actions (the top five actions in Figure 4) lack sufficient semantic information compared to other actions like keyboard inputs. For instance, *click (333, 444)* is far more abstract than *press key: enter*. To support downstream action semantic completion, PC Tracker records additional contextual information for click-related actions. Specifically, we implemented a low-overhead function `get_element_info_at_position(x, y)` that efficiently retrieves the element’s information from coordinates via system API, including its bounding box coordinates and element name (see Figure 6). When a click-related action occurs, PC Tracker passes coordinates to this function and records the clicked element’s information.

4.1.2 Observation Capture

We adopt a carefully designed timing strategy to record corresponding screenshots as state observations for actions. While it may seem intuitive to record the screenshot at the moment an action is detected, this can lead to inaccurate observation records. Figure 7 shows that humans naturally observe the environment state before making decisions and taking actions. Therefore, PC Tracker records the screenshot just before the action occurs. To achieve this, our system continuously takes screenshots and keeps the latest one in memory. When an action is detected, this latest cached screenshot (approximately 0.1 seconds prior) is recorded as the observation.

The design of only recording critical events offers notable advantages over complete video recording. It significantly reduces storage requirements, thus allowing us to capture screenshots at maximum resolution - a crucial feature for future models, as many interface details such as small-font text are only legible at high resolutions.

Besides, while many works require a complete accessibility tree to aid observation (Jia et al., 2024; Agashe et al., 2024), we find this impractical in real use. Crawling the accessibility tree requires a lengthy recursive process (taking tens of seconds to minutes), which forces users to wait after each action until the tree traversal completes, severely disrupting normal computer usage. Given the rapid progress in VLMs, we believe maintaining such exhaustive tree records is unnecessary.

4.2 Dual-Mode Collection

PC Tracker supports two recording modes: task oriented and non-task oriented, primarily differentiated by whether the interaction trajectories are associated with certain task descriptions. An overview is shown in Figure 8.

The task oriented mode records interaction trajectories that are explicitly associated with specific tasks, primarily designed to support easy annotation for supervised fine-tuning. It is divided into given task and free task modes. In the given task mode, PC Tracker randomly assigns tasks from a predefined task library to users, who then use the computer to complete the tasks while their interaction trajectories are recorded. In the free task mode, users can freely use the computer without predefined constraints and provide

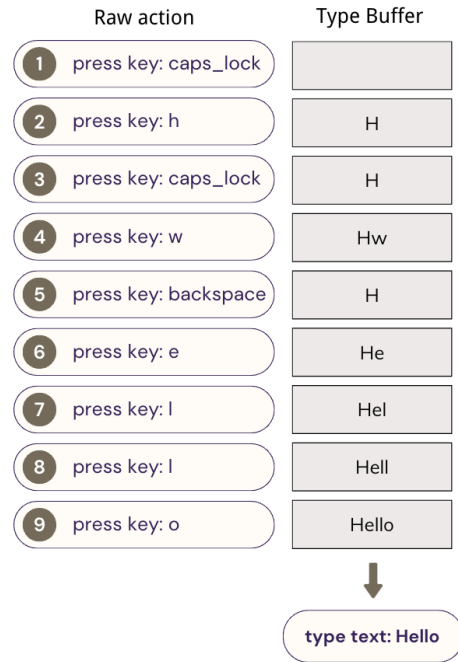


Figure 5: Example of type encapsulation.

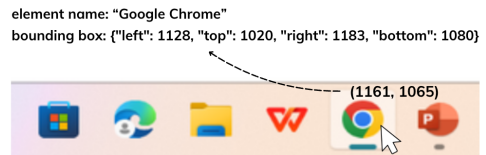


Figure 6: An example of the output from `get_element_info_at_position(x, y)` when the user clicks Chrome icon at (1161, 1065).



Figure 7: Natural flow of human interaction: Observe, Think, Act.

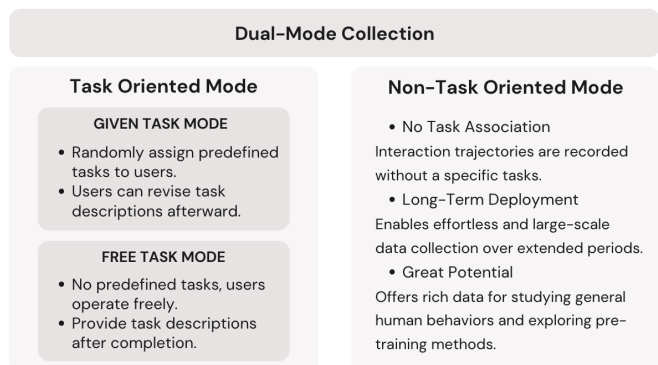


Figure 8: An overview of the dual-mode collection design

a task description for their operation sequence after completion. Notably, even in the given task mode, users have the flexibility to revise task descriptions after stopping the recording. This feature allows users to more accurately describe their activities, particularly if their operations deviated from the original assigned task or if they discovered a better way to characterize their work.

In the non-task oriented mode, the system directly records interaction trajectories without a specific task, enabling long-term deployment and effortless collection of large-scale data. While this mode may require additional approaches to infer user intentions, it provides a rich foundation for studying general human operational behaviors and exploring pre-train methods.

4.3 Privacy and Data Transparency

PC Tracker integrates privacy considerations and data transparency in its design. During recording sessions, users have full control over their data - if they notice sensitive information has been captured or if they are not satisfied with their task completion process, they can choose to discard the recording immediately without saving it. For saved trajectories, PC Tracker stores recordings locally and visualizes them in Markdown format, allowing users to easily review and check their recorded operations at any time.

5 Cognition Completion: From Raw Interaction Data to Cognitive Trajectory

Converting raw human-computer interaction data into meaningful cognitive trajectories is essential for AI to learn from human demonstrations. To achieve this goal, we conduct two-stage sequential post-processing: data refinement to ensure quality, followed by cognition completion to extract atomic action semantics and reconstruct the underlying thought process for each action.

5.1 Data Refinement

To optimize the quality of raw interaction data, we designed and implemented a comprehensive pipeline that consists of three steps: 1) trajectory filtering, 2) action filtering, and 3) standardization.

Trajectory filtering First, we examined the completeness of trajectories and files to eliminate error data caused by unexpected incidents (such as forced termination of PC Tracker), thereby ensuring robustness.

Action filtering Next, we focused on identifying and removing two categories of actions: (1) actions associated with PC Tracker, such as clicks on the start record button; (2) redundant actions in actual human operations, like control key presses for hotkey prefixes, consecutive waits due to operation pauses, and meaningless clicks. We found that these actions not only introduce unnecessary burdens to model learning but also interfere with the subsequent cognition completion process.

Standardization Finally, to ensure consistency across different data sources, we standardized all screenshot resolutions to 1080p (1920×1080), providing a standardized image input format for training.

5.2 Cognition Completion

Building upon the refined trajectories, the cognition completion stage tackles a more crucial challenge: reconstructing human cognition of computer use. This stage consists of two steps: the first step focuses on reconstructing the atomic semantics of actions, and the second step reconstructs the thought processes behind actions based on the enriched action semantics. A detailed view is shown in Figure 9.

5.2.1 Action Semantic Completion

As mentioned in Section 4.1.1, click-related actions lack semantic information due to their coordinate-based nature. To comprehensively reconstruct human cognition, we supplement click-related actions with semantic information before thought completion by generating high-quality descriptions for click targets.

PC Tracker records additional information for click elements, including the bounding box coordinates and element name (as shown in Figure 6). We highlight the click position and bounding box with red marks in the corresponding screenshot, which visualizes click-related actions on the screenshot to facilitate model understanding. This marked screenshot serves as visual input in both action semantic completion and thought completion phases.

When generating descriptions for click targets, we leverage GPT-4o to take the marked screenshot and the element name as input and generate high-quality descriptions. Notably, while PC Tracker records element names through system API, these names are often inadequate: many are missing, overly generic (e.g., numerous elements in Chrome simply labeled as “Chrome”), excessively verbose (e.g., containing complete text content for text elements), or even incorrect. In contrast, GPT-4o generates accurate and concise descriptions of click targets, facilitating the understanding of actions during thought completion.

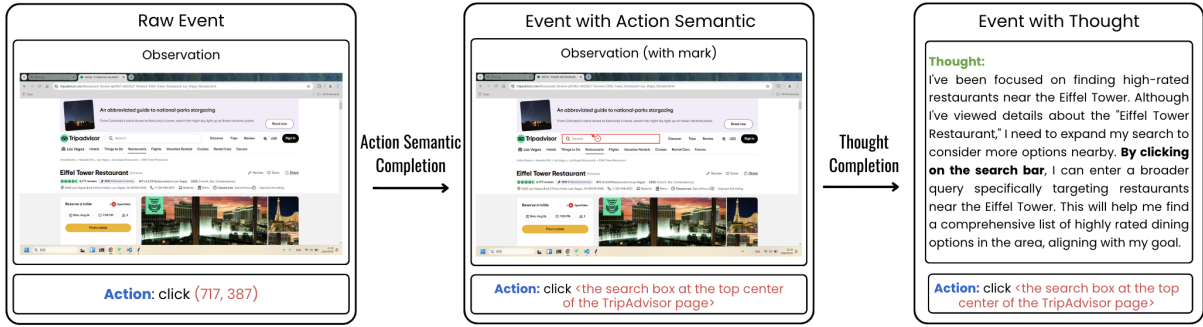


Figure 9: Visualization of our cognition completion process for a click action. (Left) Raw click event recorded by PC Tracker. (Center) Action semantic completion converts coordinates (717, 387) into a semantic description “the search box at the top center of the TripAdvisor page”. (Right) Thought completion reconstructs the human intention behind this action - finding high-rated restaurants near the Eiffel Tower by broadening the search scope.

5.2.2 Thought Completion

After completing click element information, the thought completion step reconstructs the implicit reasoning behind human actions using an iterative generation approach based on previously obtained information. Specifically, for each action in the trajectory, we provide GPT-4o with multi-dimensional input information: the task description, preceding actions with their associated thoughts, subsequent actions after semantic completion, and the corresponding marked screenshot mentioned in Section 5.2.1.

We found this comprehensive temporal context spanning past, present, and future, combined with enhanced visual information, is crucial for VLMs to understand and reconstruct human decision-making processes. Based on all this information, GPT-4o can effectively complete the accurate thought process underlying human actions, including task progression awareness, environmental exploration behaviors, and error recovery strategies.

6 PC Agent: Cognitive Agent for Complex Computer Work

This section presents our implementation of PC Agent, an AI system capable of complex computer work. Notably, our system is **built entirely on open-source models**, avoiding dependence on frontier proprietary models like GPT-4o or Claude-3.5-Sonnet.

6.1 Multi-Agent System

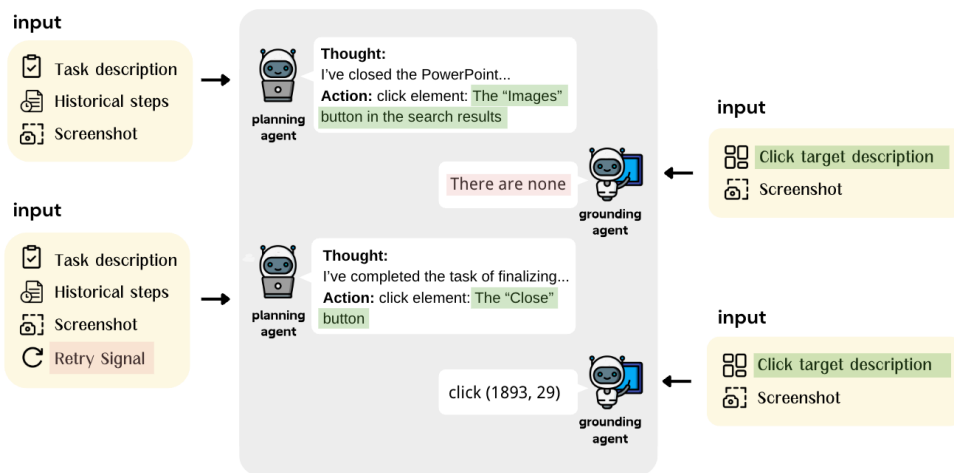


Figure 10: Illustration of multi-agent workflow. The planning agent initially attempts to click a nonexistent element *The ‘Images’ button*, which is reported by the grounding agent. Upon receiving this feedback, the planning agent reformulates its plan, and the grounding agent generates coordinates of the new click target. The workflow illustrates the error correction mechanism between agents.

In Section 2.2, we identified two major challenges of current digital agents: visual grounding and cognitive

understanding. In our preliminary implementation of PC Agent, we address these challenges through a multi-agent architecture. First, we trained a planning agent by learning from human cognitive trajectories, enabling it to acquire human cognition for effective planning. Second, we implemented a robust grounding agent with a self-validation mechanism that achieves near-human perfection. The two agents work collaboratively: the planning agent handles action decision-making, while the grounding agent executes click-related actions.

The workflow proceeds as follows, illustrated with an example in Figure 10: The planning agent first analyzes the task and observes the current state to generate a thought process and an action decision. If the action is not click-related, it will be directly executed. Otherwise, the generated target description of the click-related action will be forwarded to the grounding agent. The grounding agent not only generates the specific position coordinates of the click target but also validates grounding accuracy. If the grounding agent discovers that the planning agent is attempting to click on a non-existent target on the screen, the planning agent is prompted to reformulate its action plan. All actions are executed using the `PyAutoGUI` library.

We believe that at the current stage, the multi-agent architecture can better leverage existing VLMs’ capabilities through clear task division and is likely to achieve better results in the short term. However, looking ahead, end-to-end systems like the new Claude-3.5-Sonnet may represent the future direction of this field.

6.2 Training Planning Agent via Cognitive Trajectories

Section 2.2 posits that the limitations of current digital agents primarily stem from the lack of essential agent training data. To tackle this challenge, we utilize PC Tracker and cognition completion post-processing to efficiently build high-quality cognitive trajectories, which are then used to train the planning agent.

In the resulting training data (as shown in Figure 11), system prompts, task descriptions, and historical steps are used as textual inputs, while the corresponding screenshot (without mark) is used as the visual input. Note that each historical step consists of the thought process and action decision, incorporating the progress state within the overall task and any exploration behaviors from prior steps in natural language. For example, it allows the model to determine whether the desktop indicates the task has just started, or has been completed with all applications closed.

Furthermore, we structure the responses into distinct thought and action components. As discussed earlier, the thought process captures human cognition of computer use. This data organization enables the model to learn the underlying cognitive patterns behind actions and aligns with the ReAct (Yao et al., 2023b) paradigm where models generate verbal reasoning traces and actions in an interleaved manner.

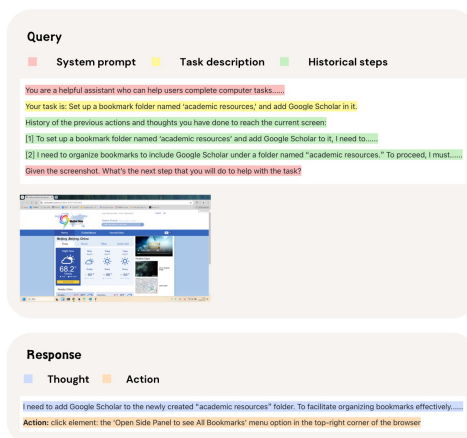


Figure 11: Training data example showing query and response structure.

6.3 Robust Visual Grounding with Self-Validation

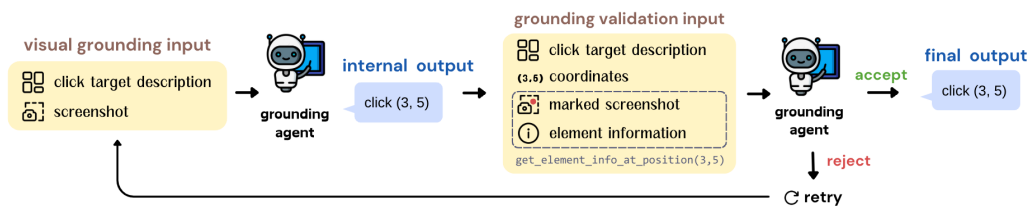


Figure 12: Illustration of the grounding agent’s self-validation mechanism.

As outlined in Section 2.2, visual grounding represents a fundamental capability for digital agents. While our exploration reveals that Molmo (Deitke et al., 2024) demonstrates remarkable proficiency in visual grounding, it still occasionally exhibits errors, particularly in the GUI domain where elements are densely arranged. However, even minor misalignments in visual grounding can lead to catastrophic failures — such as navigating to incorrect pages or triggering unintended actions — errors that current agents find particularly challenging to rectify. Our key insight is that we can leverage Molmo’s general capabilities in conjunction with external feedback from system API to implement a self-validation mechanism, further enhancing the robustness of visual grounding to achieve near-human perfection.

Specifically, as shown in Figure 12, when the grounding agent receives a click target description, it first attempts to

generate the coordinates of the target in the screenshot. This attempt leads to one of two outcomes: either the agent determines that the target does not exist and outputs “there are none”, or it successfully outputs the coordinates and retrieves the corresponding element information using the function `get_element_info_at_position(x, y)` mentioned in Section 4.1.1, including bounding box coordinates and element name. It then annotates the screenshot with red marks (similar to the marked screenshots used in cognition completion) indicating the element boundary and click position as a visual prompt, combining this with the retrieved element name to determine whether the generated position matches the initial target description. If it is judged to be inconsistent, the process reattempts output until the judgment passes or the retry limit is reached.

7 Experiment

7.1 Experimental Setup

To validate the effectiveness of our approach, we preliminarily conducted experiments on a relatively small scale, selecting the task of creating PowerPoint presentations as our primary testing ground. This choice is significant for mainly two reasons: First, presentation creation is a practical and common office task where automation can significantly reduce human workload. Second, it requires sophisticated cross-software coordination, particularly when creating high-quality presentations that involve web browsing to search and collect relevant materials. In future experiments, we’ll dive into a broader range of real-world scenarios where completing a work demands substantial human cognition, to further validate the scalability and effectiveness of our approach.

Data Collection Using PC Tracker’s task oriented mode, we collected a dataset of 133 interaction trajectories, consisting of two categories: free-form Chrome and PowerPoint related tasks (30%, 58 trajectories with an average of 29 events each) and guided PowerPoint presentation creation tasks (70%, 75 trajectories with an average of 50 events each). It is worth noting that the average action number per task in our dataset is significantly larger than those in existing datasets, as shown in Figure 13. The free-form trajectories include diverse but relatively simple operations such as web browsing, Chrome settings adjustment, and PowerPoint layout editing. In contrast, the guided presentation creation tasks follow specific presentation objectives, requiring users to engage in a more complex workflow that involves switching between PowerPoint and Chrome to collect online resources for presentation content, authentically simulating real-world work patterns.

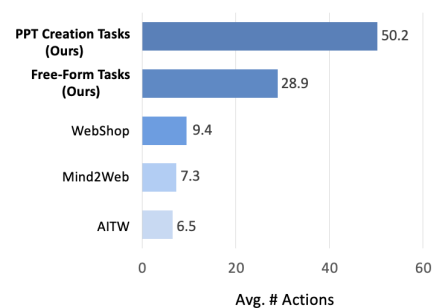


Figure 13: Average actions per task in our dataset vs. and other existing ones.

Training We fine-tuned the Qwen2-VL-72B-Instruct model with our collected dataset. The training was conducted on 32 H100 GPUs for about 2 hours. We set the context length to 8,192 tokens.

7.2 Evaluation

7.2.1 Why Not Existing Benchmarks

While benchmarks like OSWorld (Xie et al., 2024) play valuable roles in evaluating AI agents’ basic capabilities of computer use, they are inappropriate for our specific evaluation needs. Existing benchmarks tend to focus on basic operations, like duplicating the last two slides in PowerPoint. However, creating high-quality presentations is significantly more complex. The deliverables of such tasks are highly variable, often requiring subjective assessments based on aesthetics, clarity, and alignment with human preferences, which largely extends beyond the scope of existing benchmarks. Furthermore, given that our agent is currently specialized in PowerPoint presentation creation, utilizing general benchmarks to assess its performance would be both insufficient and inappropriate.

7.2.2 Human Evaluation and Case Studies

Based on these considerations, we employed human evaluation to assess PC Agent’s performance in real-world application scenarios. In our target domain — creating high-quality presentations — PC Agent exhibited outstanding performance. Its capability to execute dozens of steps and the quality of generated presentations substantially surpass existing solutions, as shown in Figure 14. We examined two representative use cases:

Case1: Complex Presentation Design We challenged PC Agent to create sophisticated presentations under detailed task specifications. It successfully accomplished the following objectives: 1) switching between PowerPoint and Chrome for resource collection, 2) organizing content into presentations with pictures, and 3) saving completed presentations to the desktop and closing all applications. Through human evaluation, we found that our

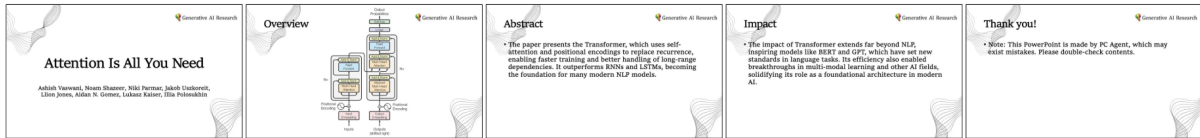


Figure 14: An example presentation created by PC Agent, see task description in Appendix C.

agent successfully handles sequences of up to 50 steps to produce multi-page presentations with practical utility, showcasing its substantial potential for real-world applications.

Case2: Batch Processing We evaluated PC Agent’s capability to handle repetitive tasks, such as creating multiple slides with consistent styling but varying content — a representative scenario where AI can significantly reduce human workload. Using a predetermined PowerPoint poster template, we tasked PC Agent with creating introduction posters for Turing Award laureates, a process requiring approximately 20 correct steps, including web-based image searches. Our assessment revealed significant potential in handling repetitive tasks: In an attempt to create 20 posters, PC Agent completed 11 posters meeting the requirements, achieving a 55% success rate. It is important to note that in some failed cases, PC Agent did create a poster, but it had subtle deviations from human requirements. This is a situation that existing benchmarks find hard to assess.

7.2.3 Analysis

Failed Case Study Our analysis revealed that mistakes stem from various ways during execution. However, these errors share a common challenge: the agent’s limited ability to recover once mistakes occur. Although our multi-agent workflow can prevent some planning mistakes through the grounding agent’s feedback, incorrect plans that get executed often lead to unrecoverable failures, indicating the need for better error recovery capabilities.

Impact of Task Description During our evaluation, we observed that the level of detail in task descriptions significantly influences PC Agent’s performance, as exemplified in Figure 15. Comprehensive task descriptions enable it to better understand human expectations, leading to more satisfactory outcomes. In contrast, brief task descriptions often result in insufficient understanding of requirements, causing the agent to favor quick but superficial solutions. Our findings suggest that optimizing task description quality is crucial for improving agent performance.

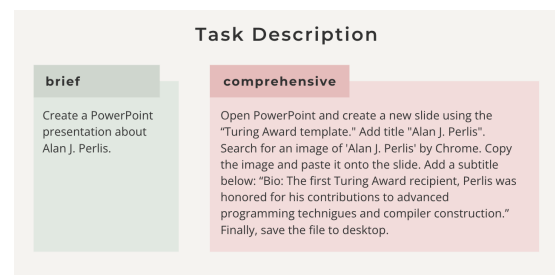


Figure 15: The brief and comprehensive task description for creating a poster about Alan J. Perlis.

Generalization Capability Our evaluation suggests that PC Agent possesses adequate performance on web browsing tasks, despite the relatively limited proportion of Chrome-related training data. This observation suggests that the operational skills acquired in specific scenarios may contribute to improved adaptability across a broader range of applications.

Data Efficiency With a relatively small training dataset of just 133 trajectories, PC Agent demonstrates remarkable task processing capabilities, particularly in PowerPoint presentation creation tasks. This data efficiency primarily stems from the high quality of the human cognitive trajectories we constructed, demonstrating the effectiveness of our framework.

8 Conclusion and Outlook

In conclusion, we presented a cognition transfer framework that efficiently guides AI to the digital world through three key components: PC Tracker for collecting human-computer interaction data, a two-stage post-processing for cognition completion, and a multi-agent system for computer task automation. Looking forward, we identify several potential directions:

Scaling and Generalization While our approach demonstrates effectiveness in presentation creation with limited training data, validating its generalization ability requires large-scale experimentation across diverse software applications and task domains.

Robustness and Long-term Planning Recent advances in LLM reasoning, exemplified by OpenAI o1 (OpenAI, 2024d), demonstrate the potential for robust operation over extended sequences. These developments suggest new approaches to long-term planning and error recovery in computer tasks.

Action Optimization Complex mouse actions, particularly **drag** and **scroll**, present notable challenges. Drag actions are especially difficult for visual grounding because both the start and end points are often abstract positions. For instance, in resizing a rectangle, the end point is an abstract position that requires an understanding of appropriate layout. Meanwhile, scroll actions are limited by our current design choice of omitting cursor movement tracking for simplification. These challenges highlight opportunities to enhance spatial understanding for abstract drag actions and reevaluate the trade-off between simplicity and completeness in cursor tracking.

Task-free Data Utilization PC Tracker’s non-task oriented mode enables large-scale collection of natural human-computer interaction trajectories. These trajectories contain rich operation patterns and behavioral preferences, which can be utilized as pre-training data to learn basic interaction patterns, as supervised learning examples by inferring user intentions after the fact, or as reference behaviors for training reward models in reinforcement learning.

Complex work Evaluation While existing benchmarks primarily focus on basic tasks with deterministic outcomes, complex real-world work like presentation creation and video editing requires more comprehensive evaluation frameworks. These frameworks need to assess deliverables through multiple dimensions, including human preference alignment, aesthetic quality, and overall completeness.

User-Friendly Task Specification While detailed task descriptions improve model performance, requiring users to provide extensive instructions can hinder system usability. Future work should explore methods to balance model requirements with user experience, such as inferring complete requirements from partial descriptions or enabling interactive clarification of task details.

Acknowledgments

We would like to thank Yan Ma, Ethan Chern, Xuefeng Li, Mingyan Yang and Xingyang Li for their valuable feedback on different stages of this work.

References

- [1] Saaket Agashe, Jiuzhou Han, Shuyu Gan, Jiachen Yang, Ang Li, and Xin Eric Wang. 2024. [Agent s: An open agentic framework that uses computers like a human.](#)
- [2] Anthropic. 2024. [Introducing computer use, a new claude 3.5 sonnet, and claude 3.5 haiku.](#) anthropic.com.
- [3] Yoshua Bengio, Aaron Courville, and Pascal Vincent. 2014. [Representation learning: A review and new perspectives.](#)
- [4] Rogerio Bonatti, Dan Zhao, Francesco Bonacci, Dillon Dupont, Sara Abdali, Yinheng Li, Yadong Lu, Justin Wagle, Kazuhito Koishida, Arthur Buckner, Lawrence Jang, and Zack Hui. 2024. [Windows agent arena: Evaluating multi-modal os agents at scale.](#)
- [5] Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. 2024. [Seeclick: Harnessing gui grounding for advanced visual gui agents.](#)
- [6] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris Callison-Burch, Andrew Head, Rose Hendrix, Favven Bastani, Eli VanderBilt, Nathan Lambert, Yvonne Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz, Aaron Sarnat, Byron Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon Borchardt, Dirk Groeneveld, Jen Dumas, Crystal Nam, Sophie Lebrecht, Caitlin Wittlif, Carissa Schoenick, Oscar Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi, Ross Girshick, Ali Farhadi, and Aniruddha Kembhavi. 2024. [Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models.](#)
- [7] Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023. [Mind2web: Towards a generalist agent for the web.](#)
- [8] Boyu Gou, Ruohan Wang, Boyuan Zheng, Yanan Xie, Cheng Chang, Yiheng Shu, Huan Sun, and Yu Su. 2024. [Navigating the digital world as humans do: Universal visual grounding for gui agents.](#)
- [9] Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. 2024. [A real-world webagent with planning, long context understanding, and program synthesis.](#)
- [10] Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. 2023. [Cogagent: A visual language model for gui agents.](#)
- [11] Jakub Hoscilowicz, Bartosz Maj, Bartosz Kozakiewicz, Oleksii Tymoshchuk, and Artur Janicki. 2024. [Clickagent: Enhancing ui location capabilities of autonomous agents.](#)
- [12] Siyuan Hu, Mingyu Ouyang, Difei Gao, and Mike Zheng Shou. 2024. [The dawn of gui agent: A preliminary case study with claude 3.5 computer use.](#)
- [13] Chengyou Jia, Minnan Luo, Zhuohang Dang, Qiushi Sun, Fangzhi Xu, Junlin Hu, Tianbao Xie, and Zhiyong Wu. 2024. [Agentstore: Scalable integration of heterogeneous agents as specialized generalist computer assistant.](#)
- [14] Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. 2024. [Visualwebarena: Evaluating multimodal agents on realistic visual web tasks.](#) *arXiv preprint arXiv:2401.13649*.
- [15] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. [Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation.](#)
- [16] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. [Visual instruction tuning.](#)
- [17] Junpeng Liu, Tianyue Ou, Yifan Song, Yuxiao Qu, Wai Lam, Chenyan Xiong, Wenhui Chen, Graham Neubig, and Xiang Yue. 2024. [Harnessing webpage uis for text-rich visual understanding.](#)
- [18] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021. [Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing.](#)
- [19] Yadong Lu, Jianwei Yang, Yelong Shen, and Ahmed Awadallah. 2024. [Omniparser for pure vision based gui agent.](#)
- [20] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. 2022. [Webgpt: Browser-assisted question-answering with human feedback.](#)

- [21] OpenAI. 2022. [Introducing chatgpt](https://openai.com). openai.com.
- [22] OpenAI. 2024a. [Gpt-4 technical report](https://openai.com).
- [23] OpenAI. 2024b. [Gpt-4v\(ision\) technical work and authors](https://openai.com). openai.com.
- [24] OpenAI. 2024c. [Hello gpt-4o](https://openai.com). openai.com.
- [25] OpenAI. 2024d. [Introducing openai o1](https://openai.com). openai.com.
- [26] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. [Hierarchical text-conditional image generation with clip latents](https://arxiv.org/abs/2204.06125).
- [27] Christopher Rawles, Sarah Clinckemaiellie, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Marybeth Fair, Alice Li, William Bishop, Wei Li, Folawiyo Campbell-Ajala, Daniel Toyama, Robert Berry, Divya Tyamagundlu, Timothy Lillicrap, and Oriana Riva. 2024. [Androidworld: A dynamic benchmarking environment for autonomous agents](https://arxiv.org/abs/2401.04202).
- [28] Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. 2023. [Android in the wild: A large-scale dataset for android device control](https://arxiv.org/abs/2305.14624).
- [29] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. [High-resolution image synthesis with latent diffusion models](https://arxiv.org/abs/2210.17723).
- [30] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. [Reflexion: Language agents with verbal reinforcement learning](https://arxiv.org/abs/2308.11432).
- [31] Open Interpreter Team. 2024. [Open interpreter official website](https://openinterpreter.com). openinterpreter.com.
- [32] Harsh Trivedi, Tushar Khot, Mareike Hartmann, Ruskin Manku, Vinty Dong, Edward Li, Shashank Gupta, Ashish Sabharwal, and Niranjana Balasubramanian. 2024. [Appworld: A controllable world of apps and people for benchmarking interactive coding agents](https://arxiv.org/abs/2401.04202).
- [33] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. [Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution](https://arxiv.org/abs/2401.04202).
- [34] Lilian Weng. 2023. [Llm-powered autonomous agents](https://lilianweng.github.io). lilianweng.github.io.
- [35] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. 2023. [Autogen: Enabling next-gen llm applications via multi-agent conversation](https://arxiv.org/abs/2308.11432).
- [36] Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao Yu, and Lingpeng Kong. 2024a. [Os-copilot: Towards generalist computer agents with self-improvement](https://arxiv.org/abs/2401.04202).
- [37] Zhiyong Wu, Zhenyu Wu, Fangzhi Xu, Yian Wang, Qiushi Sun, Chengyou Jia, Kanzhi Cheng, Zichen Ding, Liheng Chen, Paul Pu Liang, and Yu Qiao. 2024b. [Os-atlas: A foundation action model for generalist gui agents](https://arxiv.org/abs/2401.04202).
- [38] Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. 2024. [Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments](https://arxiv.org/abs/2401.04202).
- [39] Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. 2023. [Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v](https://arxiv.org/abs/2308.11432).
- [40] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2023a. [Webshop: Towards scalable real-world web interaction with grounded language agents](https://arxiv.org/abs/2308.11432).
- [41] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023b. [React: Synergizing reasoning and acting in language models](https://arxiv.org/abs/2308.11432).
- [42] Keen You, Haotian Zhang, Eldon Schoop, Floris Weers, Amanda Swearngin, Jeffrey Nichols, Yinfei Yang, and Zhe Gan. 2024. [Ferret-ui: Grounded mobile ui understanding with multimodal llms](https://arxiv.org/abs/2401.04202).
- [43] Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang, Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qingwei Lin, Saravan Rajmohan, Dongmei Zhang, and Qi Zhang. 2024a. [Ufo: A ui-focused agent for windows os interaction](https://arxiv.org/abs/2401.04202).

- [44] Chi Zhang, Zhao Yang, Jiakuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. 2023. [Appagent: Multimodal agents as smartphone users.](#)
- [45] Jiwen Zhang, Jihao Wu, Yihua Teng, Minghui Liao, Nuo Xu, Xiao Xiao, Zhongyu Wei, and Duyu Tang. 2024b. [Android in the zoo: Chain-of-action-thought for gui agents.](#)
- [46] Zhuosheng Zhang and Aston Zhang. 2024. [You only look at screens: Multimodal chain-of-action agents.](#)
- [47] Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. 2024. [Gpt-4v\(ision\) is a generalist web agent, if grounded.](#)
- [48] Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. 2024. [Webarena: A realistic web environment for building autonomous agents.](#)

A PC Tracker User Manual

1. Introduction

PC Tracker is a lightweight infrastructure for efficiently collecting large-scale human-computer interaction trajectories. The program runs seamlessly in the background, automatically capturing screenshots and keyboard & mouse activities. Figure 3 shows an example of the collected human-computer interaction trajectories.

2. Installation

- Ensure your operating system is Windows.
- Extract our software package to a location with sufficient disk space (recommended to have more than 3GB of available space for storing recorded data).

3. Quick Start

- (Optional) Set screen resolution to 16:9 (recommended 1920 x 1080).
- Open the extracted folder and launch `main.exe`.

4. Instructions

After launching the tracker, you can choose between **Task Oriented Mode** or **Non-Task Oriented Mode** for recording.

Task Oriented Mode

This mode is divided into two sub-modes: **Given Task** and **Free Task**.

Given Task In this mode, you will be assigned an uncompleted task each time.

- **Next Task:** Click `Next Task` to get the next task.
- **Previous Task:** Click `Previous Task` to return to the previous task.
- **Bad Task Feedback:** If you think the current task is difficult to complete, click `Bad Task` to discard it permanently. Alternatively, you can start the task and modify its description after completion based on your actual execution.
- **Start Recording:** Click `Start`, and the tracker window will automatically minimize while recording begins.
- **End Task:** After completing the task, click `Finish` to save the record. Or if the task execution fails or you don't want to record it, click `Fail`.
- **Modify Task Description:** After finishing the task, you can modify the task description based on your actual execution.

Free Task In this mode, you can freely use the computer and summarize the task description and difficulty yourself.

- **Start Recording:** Click `Start`, and the tracker window will automatically minimize while recording begins.
- **Save and Summarize This Record:** Fill in the task description, select difficulty (easy/medium/hard), and click `Save` to save the record.
- **Discard This Record:** Click `Discard` to discard the record.

Non-Task Oriented Mode

In this mode, you can freely use the computer, with similar methods to start and stop recording as described above.

5. Usage Notes

- **Does not currently support using extended screens.**
- **Does not currently support using Chinese input methods.**
- **Does not currently support using touchpads.**
- **The tracker window is fixed in fullscreen.** To support the filtering of tracker-related actions (such as clicking the Start button) in post-processing, the tracker window is fixed in fullscreen. You can reopen the tracker window by clicking to view the task description, then minimize it again, but please do not drag it to display in a non-fullscreen state.

6. Data Privacy

- After starting recording, your screenshots and keyboard & mouse operations will be automatically recorded. PC Tracker does not record any information from unopened software. If you believe the recording may infringe on your privacy, you can choose to discard the record.
- Collected data is saved in the `./events` folder (hidden by default). Each trajectory includes a Markdown file for easy visualization.

7. FAQ

Does the tracker have networking capabilities? PC Tracker is completely local, does not support networking, and will not upload your data.

What if my computer screen resolution is not 16:9? If your screen resolution is not 16:9, it will affect the subsequent unified processing of data. We recommend adjusting your screen resolution to 16:9.

How much space does the collected data occupy? The specific data size varies. Generally, even with intensive recording operations for 1 hour, it will not generate more than 1GB of data.

8. Contact

If you have any questions, please contact us at henryhe_sjtu@sjtu.edu.cn or zizi0123@sjtu.edu.cn.

B Prompts for Cognition Completion

B.1 Action Semantic Completion

The first stage of cognition completion is action semantic completion. Specifically, we first generate descriptions for the click targets, and then refine the generated descriptions. Table 1 and Table 2 present the prompts used for these two processes.

Table 1: Click Target Description Generation Prompt

Click Target Description Generation Prompt

Help me describe the target in the screenshot. The target may be a GUI element or an empty area on the screen.

You will be provided with:

1. A screenshot with a red mark quadruplet:
 - Frame: rectangular border around the target (may be inaccurate)
 - Circle: circle at the center of the target
 - Point: dot marking the exact click position
 - Arrow: pointing to the target
2. The name of the clicked target for reference. It's just for reference. If this name is "Unknown" or appears to be incorrect, just ignore it.

Description Rules:

1. Priority Order:
 - Highest: Circle, Point and Arrow
 - Second: Reference name (if reliable)
 - Lowest: Frame
2. Description Strategy:
 - A. For Clear GUI Elements:
 - Include position info ("top", "left", "center", etc.) if possible
 - Use visual information to describe the element
 - Refer to the provided element name if reliable
 - Examples:
 - ✓ "the button in the top-right corner of the window"
 - ✓ "the current tab at the top of the browser"
 - × "the red circle" (red marks doesn't belong to the original screenshot or element)
 - B. For Empty Areas or Uncertain Elements:
 - Focus on positional relationships
 - Use visual information to locate the position
 - Examples:
 - ✓ "empty area on the right side of the window"
 - ✓ "area near the bottom toolbar"
3. Prohibited:
 - No speculation about element functions
 - No uncertain terms like "seems", "appears", "probably"
 - No description of elements outside the frame

Output Format:

- For GUI elements: "{position description} + {element description}"
- For empty areas: "empty area + {position description}"

Examples:

- ✓ "the close button in the top-right corner of the window"
- ✓ "the 'Chrome' icon on the desktop"
- ✓ "the left thumbnail panel in current window"
- ✓ "the 'Images' tab below the search bar"
- ✓ "'click to add title'"
- ✓ "the button in the top-right corner of the browser" (when the reference name is not reliable and you are unsure about the element)
- × "what appears to be a settings button" (avoid speculation)

Important:

1. Carefully observe the screenshot and the red mark quadruplet. Use these visual cues to describe the element or position as accurately as possible. But ****DO NOT**** explicitly state the red marks in your description. Avoid phrases like "red arrow marking on the slide.." or "the red circle.."
2. When uncertain, prefer positional description over semantic or functional speculation. Be extraordinarily cautious to avoid hallucinations.
3. Be precise and output the description directly in an objective tone. Avoid sentences starting with "the target is", "The pointed target is", or "it appears to be".
4. Do not directly use the provided element name, create your own natural description based on visual information.

Note:

1. For the name of the clicked target for reference, it is either very precise or completely worthless. Judge its reliability based on visual information. If unreliable, ignore it and be cautious, preferably using only positional descriptions; if reliable, try to expand on its description as much as possible.
2. Special cases: for the text box in PowerPoint, the name of the clicked target is usually "click to add title" or "click to add text".
 - "'click to add title' ": for the title text box above the content text box or on the cover slide
 - "'click to add text' ": for the content text box below the title text box
 - "'click to add subtitle' ": for the subtitle text box below the title text box
 - "'the left thumbnail panel in current window' ": for the ****left slides thumbnail panel in PowerPoint****. But ****DO NOT**** abuse the use of "thumbnail" in other cases.

Table 2: Click Target Description Refinement Prompt

Click Target Description Refinement Prompt

You are provided with the following information about a mouse click on a computer screen:

1. A screenshot showing:
 - A red dot and circle marking the exact click location
 - A red arrow pointing to the click location
 - A red box outlining the general area of the clicked elementNote: While the dot, circle, and arrow are precise, the box might be less accurate
2. The exact coordinates of the mouse click
3. The element name from the accessibility tree
Note: This information might be incomplete, with many elements labeled as "unknown".
4. A pre-generated description of the click location
Types:
 - Empty area description (e.g., "empty area near the bottom toolbar")
 - Specific element description (e.g., "the start button on the left corner of the taskbar")

Your Task

Evaluate the provided description, determine if it is accurate. If not, provide the correct description. You can describe it as an empty area or a specific element. Do not mention the red marks on the screenshot.

B.1 Action Semantic Completion

```
# Critical Evaluation Points
1. **Priority of Visual Evidence**: The red markers (dot, circle, arrow) on the
screenshot show the ACTUAL click location. This is your primary source of truth.
But **DO NOT** explicitly state the red marks in your description. Avoid phrases
like "red arrow marking on the slide.." or "the red circle.."

2. **Element Name Usage**:
- Ignore if marked as "unknown"
- If available, use it to verify the description's accuracy
- If there's a conflict between the element name and the description, carefully
evaluate which is correct

3. **Empty Area vs. Specific Element Distinction**:
- True empty areas: Locations where clicks produce no effect
- False empty areas: Locations that appear empty but are part of specific
functional elements

# Evaluation Process
1. First, locate the exact click point using the red markers
2. Check if the provided element name offers any useful information
3. Determine if the location is a true empty area or part of a specific functional
element
4. Compare the given description against your findings
5. Provide your response based on the required format

# Important
- Carefully determine the wrong description. Most of the time, the provided
description is correct.
- The pre-generated description may have hallucinations. Carefully evaluate it.

Final Answer Format: (Response in English even the element name is Chinese)
Thought Process: {your thought process}
Answer: {your answer}

Your answer should be either:
- "Good" if the description is accurate
- "Wrong. Correct Description: {your description}" if the description is
inaccurate
```

B.2 Thought Completion

Based on the completed action semantics, we conduct thought process completion, prompt shown in Table 3.

Table 3: Thought Process Completion Prompt

Thought Process Completion Prompt
<p>You are a helpful PC Agent designed to complete tasks on a computer. Your goal is to recreate your thought process behind a specific action.</p> <p>You will be provided with:</p> <ol style="list-style-type: none">1. The task you are attempting to complete.2. A history of the steps you have already performed (up to 50, if any; none if it was the first action).3. Subsequent actions (none if this is the last action).4. The specific action you chose to take.5. A screenshot of the computer screen at the moment you decided to take the action6. The red marks on the screenshot:<ol style="list-style-type: none">A. For Click Actions (click, right click, double click):<ul style="list-style-type: none">- Frame: rectangular border around clicked element- Center: circle at element center- Click: point at exact click position- Arrow: pointing to clicked elementB. For Drag Actions:<ul style="list-style-type: none">- Start: red point and circle- End: red point and circle- Arrow: from start to end position <p>Explanation of actions:</p> <ol style="list-style-type: none">1. click element: <{element description}>: Click the element described by <code>{element description}</code>.2. right click element: <{element description}>: Right-click the element described by <code>{element description}</code>.3. double click element: <{element description}>: Double-click the element described by <code>{element description}</code>.4. drag from (x1, y1) to (x2, y2): Drag the mouse from the position (x1, y1) to (x2, y2).5. scroll (dx, dy): Scroll with offsets (dx for horizontal movement, dy for vertical movement).6. press key: key_content: Press the <code>key_content</code> on the keyboard.7. hotkey (key1, key2): Press the combination of <code>key1</code> and <code>key2</code>.8. type text: text_content: Type the text <code>text_content</code> on the keyboard.9. wait: Pause briefly, usually for system responses or screen updates.10. finish: Indicate the task has been completed.11. fail: Indicate the task has failed. <p>Further explanation of drag operation: drag from (x1, y1) to (x2, y2) is a combination of press the mouse at (x1, y1) and drag it to (x2, y2). It might has following purposes:</p> <ol style="list-style-type: none">1. Move/Translate - Moving an element from position (x1,y1) to (x2,y2) Common scenarios:<ul style="list-style-type: none">- Dragging a file/folder to a new location- Moving a UI element (window, widget) to a different position- Moving elements (shapes, text boxes, images) in a PowerPoint slide- Adjusting slider controls or resizing elements- Reordering items in a list or menu2. Range Selection - Selecting content within a rectangular region defined by (x1, y1) and (x2,y2) as diagonal points Common scenarios:<ul style="list-style-type: none">- Selecting multiple files/icons in a folder- Selecting text in a document. This is usually performed before copy/cut/delete/adjust text operation. After this action, the selected text will be highlighted.

B.2 Thought Completion

- Selecting cells in a spreadsheet
- Drawing selection rectangle on a canvas

Consider the following to give your thought process:

1. The current state of the screen and your last step (if any). Does current state align with your last plan? Are this action trying to fix something?
2. Based on the history steps, how far have you progressed in the whole task? And based on your subsequent actions, what is the expected outcome of this action? (**DO NOT** explicitly state the next action in your output.)
3. Based on all the information (task, observation, history, future), if this action seems not related to the task, is it possibly exploring the environment? Based on the above, recreate your thought process in a clear, natural first-person narrative.

Other requirements:

1. Be confident in your thought process. Avoid speculative or uncertain phrases like "it seems" or "this action might have been for."
2. You may reference future actions as context, but **DO NOT** explicitly state the next action in your explanation.
3. If there are red marks on the screenshot, you should use them to understand the action, but **DO NOT** explicitly state the red marks in your explanation. Avoid phrases like "I notice the red circles around..." or "the red arrow indicates...".
3. Keep your explanations **concise and short**, do not conduct meaningless analysis and emphasis.
4. Do not repeat the action after your thought process.

Here are some examples of the thought process:

- "I see the 'View' menu is successfully opened, so I can click the 'Slide Master' button on it to change the font for all slides."
- "To open the target powerpoint file, I should open the folder containing it first. So I need to click the folder icon on the left of the taskbar."
- "I want to click the close button to close the current window in my last step, but I see it is not closed yet in current screen. Maybe my click was not precise last time, so I need to click it again. I should click the close button on the right top corner of the window."
- "After save the file to desktop, I have successfully completed the task."
- "I need to modify the 5th slide, but it is not in the current screen. I should scroll down the page to find it."
- "I have insert a new text box and focus on it, so I can type the content now."
- "I have finished typing content in the text box. Now I can click anywhere outside the text box to deselect it and view the content on the slide."
- "I see the current file name is 'Untitled', so I should change it to a proper name. First I need to click the text box of the file name to focus on it."
- "I need to insert a new slide, so I can first click the left thumbnail panel in the PowerPoint window."
- "I need to insert a new slide, and I have clicked the left thumbnail panel in the PowerPoint window. Now I need to press key enter to insert a new slide."

Examples of thought processes for exploratory actions:

- "I need to save the file to the desktop, but I don't see a desktop option in the window. Maybe I should scroll down to see if there's a desktop option."
- "I want to select the save button, but I don't see a save option in the window. I guess I might find it by clicking the File button."
- "I need to open the settings menu, but I don't see an obvious settings icon on the current interface. Perhaps I should click on the three dots or three horizontal lines icon in the top right corner, as these often hide more options."
- "I want to change the document's font, but I can't find the font option on the toolbar. I might need to click on the 'Format' or 'Style' menu to see if I can find the font settings there."
- "I need to insert an image, but I don't see an obvious 'Insert' button. I guess I might need to right-click on a blank area of the document to see if there's an option to insert an image in the context menu."
- "I want to check the version information of this application, but I can't find the relevant option on the main interface. Maybe I should click on the 'Help' or 'About' menu, as version information is often found there."

- "I need to exit this full-screen program, but I don't see an exit button. I can try pressing the ESC key or moving the mouse to the top of the screen to see if a hidden menu bar appears."

- "I want to search for specific content on this webpage, but I don't see a search box. I can try using the shortcut Ctrl+F (or Command+F) to see if it brings up the in-page search function."

Additional PowerPoint Operation Tip:

- These steps are to add a new slide at the end of the presentation:

1. Click in the left thumbnail panel of the PowerPoint window.
2. Press the Enter key to insert a new slide.

- These steps are to add text in the text box:

1. Click 'click to add text'/'click to add title'/'click to add subtitle' to focus on the text box.
2. Type the content in the text box.
3. (Optional) Press the Enter key to finish.

Again, you are recreating your thought process when you made the action, so you should not include any post-event evaluation or similar phrases.

C Task Description Example

Table 4: The task description for Figure 14

Task Description Example

Open PowerPoint and create a new presentation using the "white template".

For the first slide, set the title to "Attention Is All You Need" and the subtitle to "Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin."

For the second slide, set the title to "Overview". Search for an image of "Transformer Attention Is All You Need" by Chrome and copy the source image. Back to PowerPoint and paste the image into the "Overview" slide.

For the third slide, title it "Abstract", with the following text: "The paper presents the Transformer, which uses self-attention and positional encodings to replace recurrence, enabling faster training and better handling of long-range dependencies. It outperforms RNNs and LSTMs, becoming the foundation for many modern NLP models."

For the fourth slide, title it "Impact", with the following text: "The impact of Transformer extends far beyond NLP, inspiring models like BERT and GPT, which have set new standards in language tasks. Its efficiency also enabled breakthroughs in multi-modal learning and other AI fields, solidifying its role as a foundational architecture in modern AI."

For the last slide, add title "Thank you!" and following text: "Note: This PowerPoint is made by PC Agent, which may exist mistakes. Please double-check contents."

Save the presentation and close PowerPoint and Chrome window and you're all set!
