


UFO : A UI-Focused Agent for Windows OS Interaction

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
Microsoft
UFO-Agent@microsoft.com

Abstract

We introduce UFO , a **UI-Focused** agent designed to fulfill user requests tailored to Windows OS applications by observing and analyzing the GUI and control information of these applications. UFO utilizes a hierarchical dual-agent framework that decomposes user requests using a divide-and-conquer approach, enabling seamless navigation and addressing sub-tasks across multiple applications. It also incorporates a control automator tailored for Windows OS, which detects control elements effectively and allows for fully automated execution. As a result, UFO simplifies complex and time-consuming processes into tasks that can be completed with natural language commands.

We conducted testing of UFO across 9 popular Windows applications, encompassing a variety of scenarios. The results derived from both quantitative metrics and real-case studies, underscore the superior effectiveness of UFO in fulfilling user requests. To the best of our knowledge, UFO stands as the first UI agent specifically tailored for task completion within the Windows OS. The open-source code for UFO is available on <https://github.com/microsoft/UFO>.


1 Introduction

The application of Visual Large Language Models (VLMs) has demonstrated transformative potential in addressing complex multimodal challenges, emulating human-like reasoning, planning, and collaboration (Zhang et al., 2023; Yang et al., 2023c; Wu et al., 2023; Zheng et al., 2024). Graphical User Interfaces (GUIs), though designed primarily for human interaction, offer a critical interface through which VLMs can engage with software applications. By perceiving and analyzing visual information from UI elements, VLMs can ground their actions in the real world, bridging the gap

between human-like interactions and automated processes (Hong et al., 2023). Such interactions can potentially enable VLMs to fulfill user requests expressed in natural language and execute actions on physical devices, thereby translating complex requests into tangible outcomes.

In this context, the Windows operating system (OS) emerges as a representative platform due to its frequent use in daily life (Adekotujo et al., 2020), the variety of versatile applications built upon it (Ramler et al., 2018), and the complexity of tasks that often require long-term planning and interaction across multiple applications (Stallings, 2005). Additionally, the system’s emphasis on GUI interaction highlights the need for a general intelligent agent capable of understanding user requests in natural language and autonomously interacting with the UIs of Windows-based applications.

The potential of such an agent is compelling, as it promises to simplify both routine and creative tasks, significantly reducing the user’s effort for learning and operating within various applications. However, Windows presents unique challenges due to its diverse applications, larger screen resolutions, and more complex GUIs compared to platforms like mobile or Linux. A naive agent may struggle to understand and interact with the complex interfaces, making it difficult to fulfill user requests, especially those that span multiple applications.

To address this gap, we introduce UFO , a specialized **UI-Focused** VLM Agent designed for seamless interaction with the Windows OS. We illustrate this concept at a high level in Figure 1, where UFO composes and sends an email by integrating text extracted from a Word document, observations from an image, and summaries from a PowerPoint presentation. UFO employs a dual-agent framework that utilizes a divide-and-conquer approach to fulfill user requests. In this framework, the HostAgent decomposes the user request

*Corresponding author.

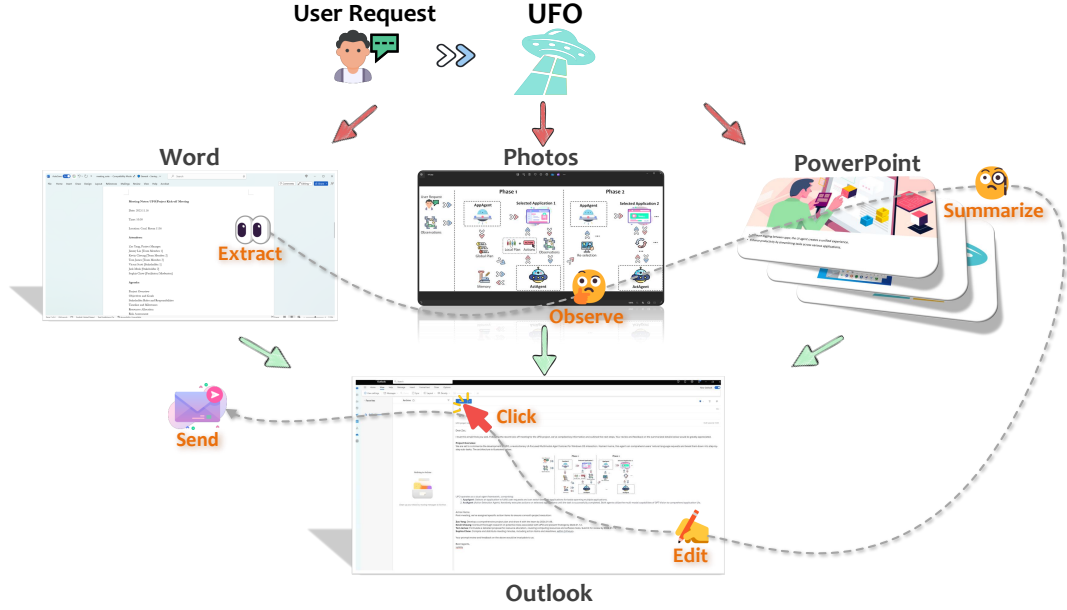


Figure 1: The high-level concept of UFO. It completes a user request by navigating through multiple applications, all accomplished **automatically**.

into sub-tasks, while the AppAgent executes the assigned sub-tasks within individual applications. Each agent utilizes screenshots and extracts information from the GUI to make informed decisions. To ground actions in the application, UFO incorporates a control automator component specifically tailored for the Windows OS. This inspects control elements and translates actions from the agents into grounded executions within applications. This setup ensures reliable action execution and complete automation without the need for human intervention.

As user requests often involve interactions across multiple applications, UFO includes a switching mechanism within the HostAgent to facilitate transitions between different applications. This capability allows UFO to handle complex, multi-application tasks that are frequently overlooked by other agents. Additionally, UFO is highly extensible, allowing users to design and customize actions for specific tasks and applications. It also incorporates a safeguard feature to prevent the execution of risky actions. These enable UFO to automate lengthy and tedious processes through natural language commands while ensuring safety.

We conducted comprehensive testing of the UFO framework on our WindowsBench benchmark which include 50 tasks involving 9 widely used Windows applications to capture a range of scenarios reflective of users’ daily needs. The evaluation included both quantitative metrics and detailed case studies, demonstrating that UFO signif-

icantly outperforms other GUI agents. This performance underscores the robustness and adaptability of our design, particularly in handling complex requests that span multiple applications. To the best of our knowledge, UFO represents the pioneering agent specifically tailored for general applications within the Windows OS.


2 Related Work

The use of multimodal large language models (LLMs) for navigating and controlling GUIs in applications has become a prominent area of research. AppAgent (Yang et al., 2023b) leverages GPT-4V to mimic smartphone users, enabling actions on mobile applications based on smartphone snapshots, thereby autonomously fulfilling user requests on physical devices. MobileAgent (Wang et al., 2024) enhances the capabilities of GPT-4V within a similar mobile agent framework by integrating Optical Character Recognition (OCR) tools, allowing it to achieve task completion rates comparable to human performance. In contrast, OS-Copilot extends this concept to desktop environments, using self-directed learning to continuously evolve the agent’s capabilities. Cradle applies DINO (Liu et al., 2023) and SAM (Kirillov et al., 2023) for control detection within applications, further extending these techniques into the realm of video game playing (Tan et al., 2024).

Distinguishing itself from existing frameworks, UFO stands out through its unique features: (i) it is tailored to the Windows OS with specialized

control inspection and interaction; (ii) it employs a dual-agent, divide-and-conquer approach for handling cross-application tasks; and (iii) it incorporates several design elements to ensure both customization and safety.

3 The Design of UFO

We present UFO , a **UI-Focused Agent** designed for the Windows OS interaction. UFO comprehends users' requests expressed in natural language, and observe the UI screenshots of applications and operates on their control elements to fulfill the overall objective, transcending the boundaries of different applications.

3.1 UFO in a Nutshell

First, we present the workflow of UFO in Figure 2. UFO operates as a centralized dual-agent framework consisting of (i) the **HostAgent**, which decomposes user requests into sub-tasks and assigns them to the appropriate AppAgent in a divide-and-conquer manner, and (ii) the **AppAgent**, responsible for iteratively executing actions on a local application until the task is successfully completed within its scope. Both agents leverage the multi-modal capabilities of VLM to understand the application UI and fulfill the user's request. They also utilize a **Control Automator** module tailored to Windows OS to detect controls and ground their actions, enabling tangible impacts.

Upon receiving a user request, the HostAgent first analyzes the intent and formulates a global plan consisting of sub-tasks, each to be completed within a specific application. The AppAgent then accepts its assigned sub-task and carries it out within the designated application by interacting with the application's controls via a control automator module. If the user request spans multiple applications, the HostAgent generates multiple sub-tasks, assigning them to the AppAgent to complete sequentially. This iterative process continues until the entire request is fully executed.

3.2 HostAgent

The HostAgent is responsible for analyzing user intent and decomposing the request into sub-tasks, each of which operates on a single application. We show its functionality in Figure 3.

Observation. The HostAgent gathers the following inputs from the user and the system:

- **User Request:** The original query submitted by the user to UFO.
- **Desktop Screenshots:** Images capturing the current state of the desktop.
- **Application Information:** Active applications, including their names and types.

This information provides the HostAgent with a comprehensive data to support decision-making. Desktop screenshots and application details enable the HostAgent to understand the current environment and set boundaries for task decomposition and application selection. Additionally, the HostAgent is equipped with several demonstrated examples to activate in-context learning (ICL) (Brown et al., 2020).

Prediction. After collecting all necessary information, the HostAgent uses GPT-V to generate the following outputs:

- **Observation:** A detailed description of the current desktop based on the screenshots.
- **Thoughts:** The logical thinking required to complete the task, following the Chain-of-Thought (CoT) paradigm (Wei et al., 2022).
- **Selected Application:** The chosen application for the next sub-task.
- **Status:** The overall task status, denoted as either "CONTINUE" or "FINISH".
- **Plan:** A plan for sub-task decomposition.
- **Comment:** Additional remarks or information for the user, including a brief progress summary and key points to be noted.

By prompting the HostAgent to provide its observations and thoughts, we encourage it to analyze the current status, offering a clear explanation of its logic and decision-making process. This approach enhances the logical coherence of its decisions (Wei et al., 2022; Ding et al., 2023) and contributes to overall interpretability. The HostAgent may also leave comments for the user, highlighting potential issues or addressing any queries. Once the HostAgent decompose the overall goal, UFO proceeds to execute each sub-task through the AppAgent within that application.

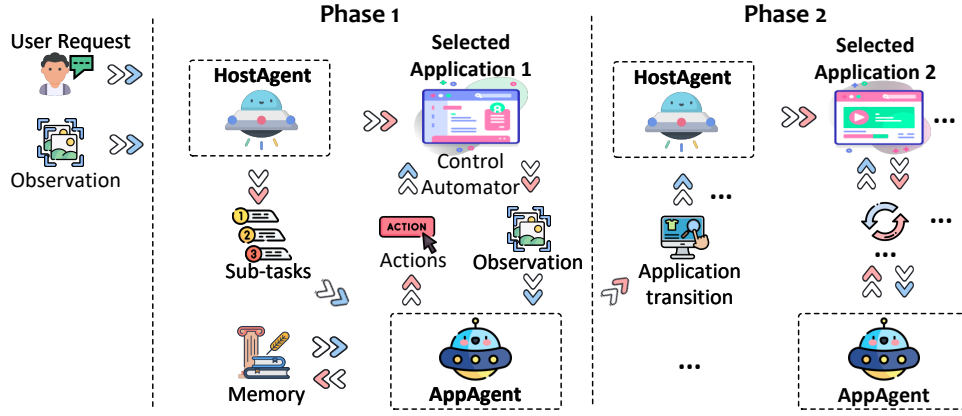


Figure 2: The overall workflow of the UFO.

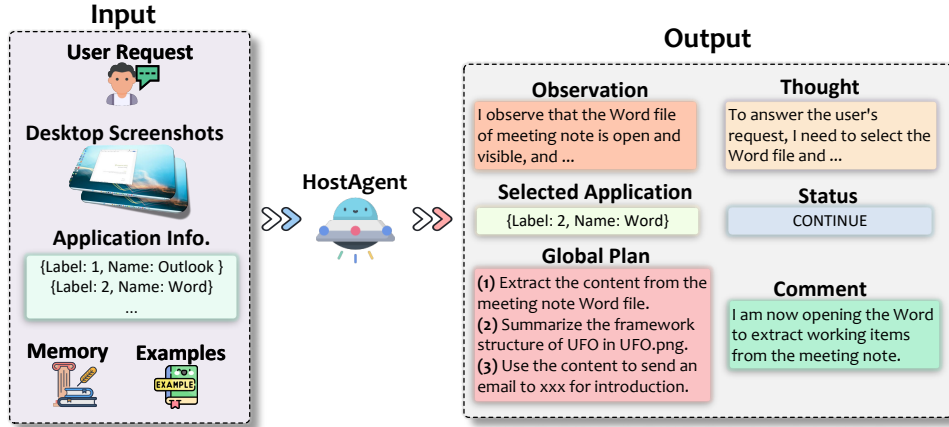


Figure 3: An illustration of the HostAgent.

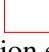
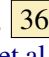
3.3 AppAgent

The AppAgent functions as the downstream entity responsible for completing the sub-task assigned by the HostAgent on the selected application, as illustrated in Figure 4.

Observation. At each operation step, AppAgent gathers the following input:

- **User Request:** The original query submitted by the user to UFO.
- **Sub-Task:** The specific sub-task assigned by the HostAgent.
- **Screenshots:** The application screenshots are provided in three forms: (i) Previous Screenshot; (ii) Clean Screenshot; and (iii) Annotated Screenshot.
- **Control Information:** A list of control names and types that are enabled for operations within the selected application.

All of these information are gathered automatically. Unlike the HostAgent, UFO provides the AppAgent with three types of screenshots to enhance its decision-making process. The previous screenshot, which highlights the last selected control in

a red rectangle (*i.e.*, ) helps the AppAgent understand the operation executed in the last step. The clean screenshot offers a clear view of the application’s current status without any annotation obstructions, while the annotated screenshot labels each control with a number (*e.g.*, ) using the Set-of-Mark (SoM) method (Yang et al., 2023a), facilitating a better understanding of the UI elements’ functions and locations. Additionally, the AppAgent employs demonstrated examples for ICL, refining its capabilities.

Prediction. The AppAgent analyzes all the inputs and outputs the following:

- **Observation:** A detailed description of the current application window.
- **Thoughts:** The logical reasoning and rationale behind the current action decision.
- **Selected Control:** The label and name of the chosen control for the operation.
- **Function:** The specific function and its arguments that applied to the selected control.

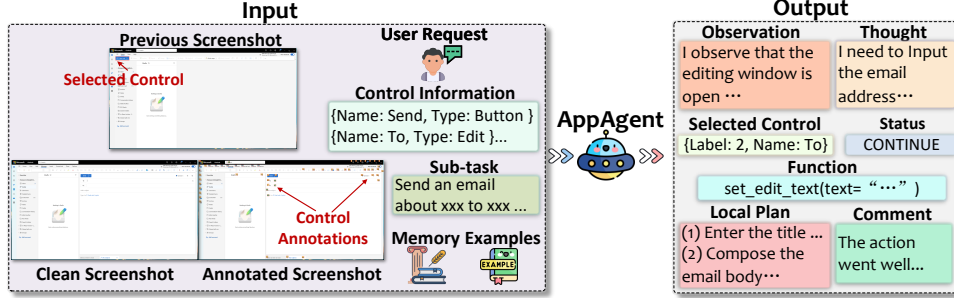


Figure 4: An illustration of the AppAgent.

- **Status:** The task status, indicated as one of the following: “CONTINUE” if further action is needed, “FINISH” if the task is complete, “PENDING” if user confirmation is required, or “APP_SELECTION” when the task is complete on the current application and a switch to another application.
- **Plan:** A more precise and fine-grained plan for future actions on the application.
- **Comment:** Additional comments or information, such as a brief progress summary, key points to note, or changes in the plan.

The AppAgent relies on its observations and CoT reasoning to carefully evaluate its decisions. The function applied to the selected control is then executed through the control automator module. The plan guides future actions, while comments facilitate communication with the user.

The AppAgent continues this cycle of observing and acting step-by-step until the assigned sub-task is fully completed. Once a sub-task is finished, the HostAgent evaluates whether all sub-tasks are finished. If any sub-tasks remain, the HostAgent will assign them to another AppAgent for further processing.

3.4 Control Automator

To provide control and application information inputted to UFO and execute actions from the AppAgent within the application, UFO utilizes the Python package `pywinauto` (Bim and Minshuai, 2014), which provides essential tools for inspecting UI controls and performing automated actions on them. This package offers a variety of functions for automating GUI interactions, enabling programmatic control of Windows applications. For the backend, we selected the UI Automation (UIA) API (Dinh et al., 2018) due to its robust and universal support for UI inspection and

interaction, specifically tailored for the Windows environment.

Inspection. UFO utilizes `pywinauto` to inspect all actionable controls within an application, retrieving their precise location and bounding box to facilitate annotations using the Set-of-Mark (SoM) method (Yang et al., 2023a) (see an example in Figure 6 in the Appendix). Additionally, `pywinauto` provides rich contextual information for each control, including its title and type, which are critical for the AppAgent’s control and action selection.

The specialized control inspection using the UIA is a key factor in the success of UFO. This approach provides more reliable results compared to model-based control detection methods like DINO and SAM used by GUI agents. Leveraging the UIA framework, this transition from segmenting and predicting targeted controls to selecting from a predefined list leads to significantly higher task completion rates.

Interaction. For the specific functions applied to controls, we utilize common and widely-used mouse operations supported by `pywinauto`, along with custom-developed actions. These actions include `Click`, `SetText`, `GetText`, `Scroll`, `Annotate`, and `Summary`.

The `Click` and `Scroll` functions replicate standard mouse operations on a control or specific coordinates. `SetText` simulates typing, including the use of hotkeys, similar to keyboard input. `GetText` extracts the text from a text box. `Annotate` and `Summary` functions are custom operations designed to meet specific needs within UFO. The `Annotate` function enables the re-annotation of the GUI with a more concise list of controls (details in Section 3.6.3), while the `Summary` function allows UFO to describe its visual observations in text to fulfill user requests. A detailed description of these actions is provided in Section A.2.

At each step, the AppAgent selects an appropriate action from this list to execute on the chosen UI control. Through the control automator, UFO can make tangible and effective changes to the system.

3.5 Memory

To enhance the capabilities of UFO, both the HostAgent and AppAgent maintain a memory module that records historical information from previous steps, including plans, thoughts, comments, actions, and execution results. This information allows each agent to analyze past trajectories and minimize the likelihood of repeating ineffective actions. Moreover, the memory module facilitates cross-application communication. For example, execution results, such as text extracted from a document, are stored in the memory module. The AppAgent can then incorporate this information for subsequent actions, such as composing an email using text from results obtained in earlier steps.

3.6 Special Design in UFO

UFO incorporates a series of dedicated design elements tailored to the Windows OS, to facilitate more effective, automated, and secure interactions with UI controls, as detailed next.

3.6.1 Human in the Loop

User requests may sometimes be ambiguous, and the tasks completed by UFO might not always meet expectations. To address this, UFO provides users with the ability to engage in iterative interactions rather than relying solely on one-shot task completions. After a task is completed, users can request UFO to refine the previous task, propose a new task, or even perform operations to assist UFO in areas where it may lack proficiency, such as providing a username input. This human-in-the-loop approach not only sets UFO apart from other existing UI agents but also allows it to incorporate user feedback, facilitating continuous improvement.

3.6.2 Action Customization

UFO allows users to register their customized operations like plugins beyond the constraints of UIA. This involves specifying the purpose, arguments, return values, and providing illustrative examples for demonstration purposes. Once the registration process is completed, the customized operation becomes available for execution by UFO. This inherent flexibility renders UFO a highly extendable

framework, enabling it to fulfill more intricate and user-specific requests within the Windows OS.

3.6.3 Control Filtering

UIA can detect hundreds of control items in a application. However, annotating all these controls can clutter the application UI screenshots, obstructing the view of individual items and generating an extensive list that may pose challenges for UFO in making decisions.

To address this, UFO employs a dual-level control filtering approach, comprising the hard level and the soft level. At the hard level, candidate controls are constrained based on specific control types with high relevance and popularity, as detailed in Section A.3. Moreover, we incorporate a soft filter mechanism that empowers UFO to dynamically determine whether to re-select a more concise list of specified controls. This adaptive filtering is triggered when UFO perceives an excessive number of controls, potentially obstructing the visibility of the required one. In such scenarios, UFO captures a new screenshot annotated only with these filtered controls, facilitating a more focused and effective filtering process.

3.6.4 Plan Reflection

While both the HostAgent and AppAgent are responsible for formulating plans, the actual state of the application’s UI may not always match the anticipated conditions. To manage this dynamic environment, UFO is designed to continuously revise its initial plan at every decision step. This adaptive approach allows UFO to respond effectively to the evolving status of the application based on its ongoing observations. The effectiveness of these reflective mechanisms enhances UFO’s ability to correct errors and explore alternative solutions, as supported by recent studies on task execution and reflection (Qiao et al., 2023; Shinn et al., 2023).

3.6.5 Safeguard

Lastly, we acknowledge the sensitivity of certain actions within the system, such as irreversible changes resulting from operations like file deletion. To reduce these risks, UFO incorporates a safeguard mechanism to seek user confirmation before executing such actions. The safeguard functionality is not limited to the following enumerated list in Table 5 in Appendix, as UFO intelligently assesses the sensitivity of each action. Users can also define their own sensitive actions, which require confirmation before execution. With the deployment of this

safeguard, UFO establishes itself as a significantly safer and trustworthy agent, mitigating the risk of compromising the system or jeopardizing user files and privacy. This feature is not addressed by other GUI agents.

4 Experiment

4.1 Experiment Setting

Benchmark. We developed a benchmark called WindowsBench, designed by five expert Windows users. This benchmark consists of 50 user requests across nine widely-used Windows applications that are essential for daily tasks. The selected applications include Outlook, Photos, PowerPoint, Word, Adobe Acrobat, File Explorer, Visual Studio Code, WeChat, and Edge Browser. These applications serve various purposes, ensuring a diverse and comprehensive evaluation.

For each application, we created five distinct requests, with an additional five requests involving interactions that span multiple applications. This setup provides a total of 50 requests, where each application includes at least one request linked to a subsequent task, allowing for a thorough assessment of UFO’s interactive capabilities. A detailed list of the requests used in WindowsBench can be found in Table 6, 7, and 8 in Appendix Section C.

Baselines. We consider two sets of baselines: non-agentic and agentic methods. For non-agentic methods, we selected GPT-3.5 and GPT-4, instructing them to provide step-by-step instructions to complete the user requests. Since these models cannot directly interact with applications, a human surrogate executes the operations according to their guidance. For agentic frameworks, we chose open-source frameworks, OS-Copilot (Wu et al., 2024) and Cradle (Tan et al., 2024), all of which can operate on computers autonomously. These agents primarily rely on model-based methods for control inspection. All agents, including UFO use GPT-4V as the inference engines. Detailed settings are provided in Section B.

Metrics. We assess UFO from 4 perspectives: success, step, completion rate, and safeguard rate. The success metric determines if the agent successfully completes the request. The step refers to the number of actions the agent takes to fulfill a task, serving as an indicator of efficiency. The completion rate is the ratio of the number of correct steps to the total number of steps. Lastly, the safeguard

Table 1: Comparison on WindowsBench.

Framework	Success	Step	Completion Rate	Safeguard Rate
GPT-3.5 (Human Surrogate)	24%	7.86	31.6%	50%
GPT-4 (Human Surrogate)	42%	8.44	47.8%	57.1%
OS-Copilot	58%	7.14	62.3%	0%
Cradle	70%	6.33	75.9%	0%
UFO	86%	5.48	89.6%	85.7%

Table 2: The detailed performance breakdown across applications achieved by UFO.

Application	Success	Step	Completion Rate	Safeguard Rate
Outlook	100.0%	6.8	94.0%	100.0%
Photos	80.0%	4.0	96.7%	100.0%
PowerPoint	80.0%	5.6	88.8%	50.0%
Word	100.0%	5.4	92.7%	-
Adobe Acrobat	60.0%	4.2	78.7%	100.0%
File Explorer	100.0%	4.8	88.7%	100.0%
Visual Studio Code	80.0%	4.0	84.0%	-
WeChat	100.0%	5.0	98.0%	66.7%
Edge Browser	80.0%	5.2	92.0%	100.0%
Cross-Application	80.0%	9.8	83.0%	100.0%

rate measures how often UFO requests user confirmation when the request involves sensitive actions. All metrics are assessed through human evaluation. More details on the evaluation are presented in Appendix Section C.

4.2 Performance Evaluation

Non-Agentic Comparison. We first compare UFO with non-agentic models, as presented in Table 1. Notably, UFO achieves an 86% success rate, more than double that of GPT-4, highlighting its advanced capability in executing tasks on the Windows OS. Furthermore, UFO demonstrates the highest completion rate, indicating its precision in taking accurate actions. It also completes tasks with the fewest steps, underscoring its efficiency, whereas GPT-3.5 and GPT-4 tend to require more steps but are less effective overall.

The inferior performance of the non-agentic methods can be attributed to two key factors. First, these baselines lack the ability to directly interact with the application environment, relying on human surrogates for execution. This reliance hinders their ability to adapt to changes and dynamic reflections in the environment. Second, these models process only textual input, overlooking the importance of visual information in GUI interactions, which is often crucial for effective task completion.

Agentic Comparison. Turning to the agentic approaches, UFO remains the superior framework, significantly outperforming the best baseline, Cradle, by 16% in success rate. This advantage is primarily due to two factors. First, baseline frame-

Table 3: Ablation study of different setting of UFO.

Framework	Success	Step	Completion Rate	Safeguard Rate
UFO w/o screenshots	72%	7.84	74.2%	82.1%
UFO w/o CoT	76%	5.61	79.7%	78.3%
UFO w/o SoM	80%	6.21	84.6%	83.0%
UFO w/o Control Filter	78%	6.57	81.52%	84.8%
UFO Qwen-VL-Plus	64%	7.90	64.2%	73.2%
UFO Genimi 1.5 Pro	78%	6.07	83.9%	81.2%
UFO GPT-4V	86%	5.48	89.6%	85.7%

works typically rely on model-based methods like SOM and DINO for control element detection. However, these models often fail when faced with densely packed control layouts, leading to unsuccessful task execution. Second, while some base-lines employ a HostAgent-like manager to decompose complex tasks across multiple applications, they struggle with cross-application tasks. In contrast, UFO leverages the UIA framework, specifically tailored for Windows, for robust control detection. Its divide-and-conquer strategy effectively handles cross-application tasks, making it more adept at navigating the Windows OS.

The safeguard feature is not included in the baseline methods, resulting in a safeguard rate of 0%. In contrast, UFO achieves a safeguard rate of 85.7%, demonstrating its ability to accurately classify sensitive requests. This highlights UFO’s superior security and reliability.

App-Breakdown. We present the breakdown of UFO’s performance across different applications in Table 2. The “-” in the safeguard rate column denotes that all requests are not sensitive. Notably, UFO demonstrates strong performance across all applications, showcasing its effectiveness in operating on versatile software. An exception is Adobe Acrobat, where UFO achieves a 60% success rate. This is because many control types in Adobe Acrobat are not supported by UIA, posing challenges for UFO in operation.

Importantly, when cross-application requests (last row), UFO maintains a high level of performance. Despite requiring more steps (average of 9.8) to fulfill, UFO achieves an 80% success rate, an 83% completion rate, and a 100% safeguard rate. This underscores UFO’s sophistication in navigating across different applications for complex tasks, proving the superiority of our divide-and-conquer design.

4.3 Ablation Study

We conducted an ablation study to assess the impact of excluding specific components from UFO,

including visual information (screenshots), CoT, SoM annotation, and control filtering. We also compared the performance of Qwen-VL-Plus (Bai et al., 2023) and Genimi 1.5 Pro (Reid et al., 2024) with GPT-4V.

The results reveal that omitting any of these components leads to a noticeable decline in performance compared to the full UFO. Specifically, removing visual information results in a 14% drop in the success rate, underscoring the critical role of visual data in enabling accurately understand and interact with application GUIs. The absence of CoT causes a 12% reduction in the success rate, highlighting the importance of step-by-step reasoning for handling complex user requests. Removing SoM annotation from screenshots leads to a performance decrease, demonstrating the importance of accurate control element localization for effective task execution. Also, removing control filter leads to 4% drop of success rate, proving its effectiveness for reducing redundant information. These confirm that each component plays a vital role in UFO’s overall effectiveness and efficiency. Additionally, GPT-4V consistently outperforms Qwen-VL-Plus and Genimi 1.5 Pro across all metrics, establishing it as the superior inference engine within UFO.

4.4 Case Study

We make a complex request shown in Figure 5, “My name is Zac. Please read the meeting note to identify all action items, and include the detailed description LLM training workflow in the LLM-training.png, to compose a new email of these contents. Send the full email to our leader Hidan via email address to ask for his review”. This request puts a significant demand on UFO’s ability to navigate and interact across 3 different applications.

In response, UFO formulates a global plan navigating between Word, Photos, and Outlook to achieve the overarching objective. Initially, it activates the required document in Word, utilizing the `GetText` API to extract text from the canvas, and logged into its memory. Subsequently, UFO switches to the LLM-training image file in Photos, observing and generating a detailed description for future email composition. With all the necessary information gathered, UFO opens the Outlook, accessing the “New Email” button to initiate the editing block. The agent then autonomously inputs the email recipient, drafting a subject, and composing the email body, including all required information. Prior to sending, the **safeguard** prompts the user

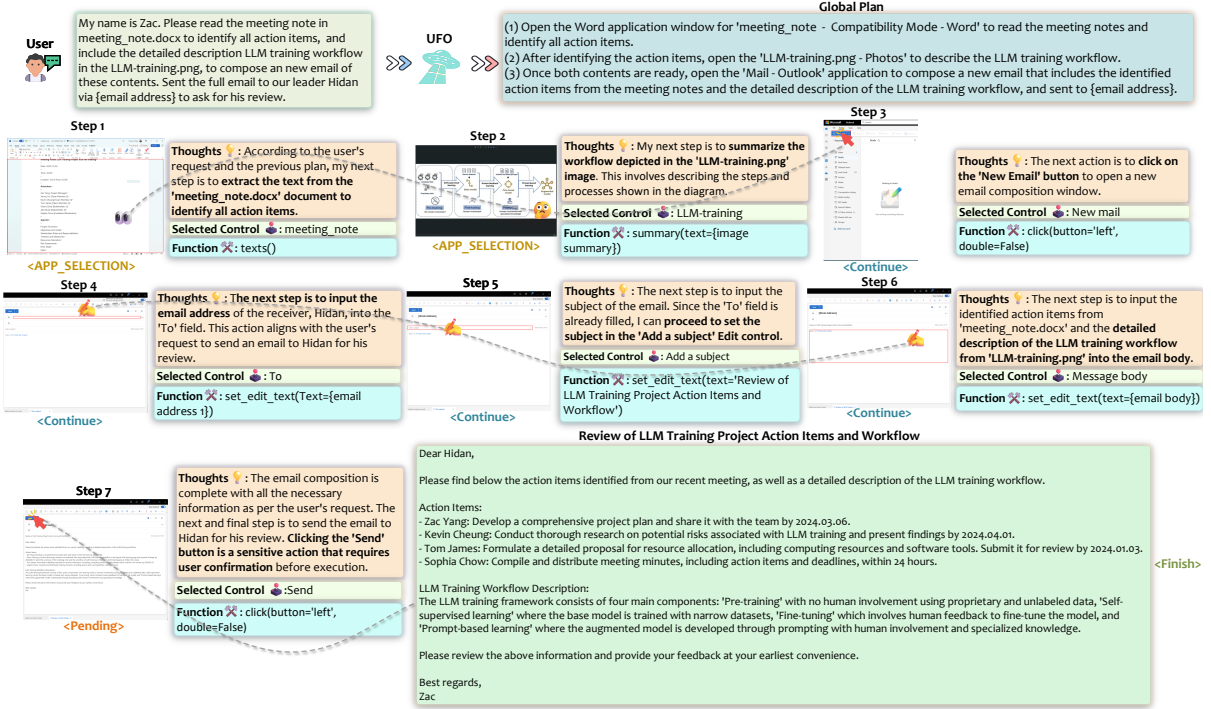



Figure 5: A example of UFO composing email with sources from multiple applications.

for confirmation due to the sensitive nature of the action. Once confirmed, the email is sent. This illustrates the remarkable ability of UFO in completing complex tasks that require long-term planning, even when spanning multiple applications.

For further detailed case studies, please refer to Appendix Section F.

5 Conclusion

We present UFO , a UI-focused agent designed to fulfill user requests in natural language by intelligently interacting with applications on the Windows OS. UFO analyzes GUI and control information to select the most appropriate application and element to execute actions that satisfy user queries. With its dual-agent architecture, UFO efficiently manages complex task switching between applications, enabling the completion of complex tasks that span multiple software environments. Key features such as action customization and safeguard further enhance UFO's extensibility and safety. Evaluation results from 50 requests across 9 popular Windows OS applications demonstrate its exceptional versatility and generalizability. To the best of our knowledge, UFO is the first UI automation agent specifically tailored for the Windows OS.

6 Limitations

The available UI controls and actions are currently limited by those supported by pywinauto and Windows UIA. To broaden UFO's capabilities, we plan to expand its scope by supporting alternative backends, such as Win32 API, or incorporating dedicated GUI grounding models for visual detection, as demonstrated by CogAgent (Hong et al., 2023). This enhancement will enable UFO to operate across a broader range of applications and handle more complex actions.

Secondly, we recognize the challenge UFO faces when exploring unfamiliar application UIs, which may be niche or uncommon. To address this, we propose leveraging knowledge from online search engines as an external knowledge base for UFO. Analyzing both textual and image-based guidelines in search results will empower UFO to distill a more precise and detailed plan for completing requests on unfamiliar applications.

Finally, when using UFO for task completion, the user's use of the mouse or keyboard may interrupt the agent's execution, potentially leading to a negative user experience due to the exclusive control required by UFO. We acknowledge this issue and will explore new paradigms of human-agent interaction to enhance the overall experience and mitigate such interruptions.

References

- Akinlolu Adekotujo, Adedoyin Odumabo, Ademola Adedokun, and Olukayode Aiyeniko. 2020. A comparative study of operating systems: Case of windows, unix, linux, mac, android and ios. *International Journal of Computer Applications*, 176(39):16–23.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond.
- HE Bim and WANG Min-shuai. 2014. Application of pywinauto in software performance test. *Computer and Modernization*, (8):135.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Yuhang Chen, Chaoyun Zhang, Minghua Ma, Yudong Liu, Ruomeng Ding, Bowen Li, Shilin He, Saravan Rajmohan, Qingwei Lin, and Dongmei Zhang. 2023. Imdiffusion: Imputed diffusion models for multivariate time series anomaly detection. *arXiv preprint arXiv:2307.00754*.
- Ruomeng Ding, Chaoyun Zhang, Lu Wang, Yong Xu, Minghua Ma, Wei Zhang, Si Qin, Saravan Rajmohan, Qingwei Lin, and Dongmei Zhang. 2023. Everything of thoughts: Defying the law of penrose triangle for thought generation. *arXiv preprint arXiv:2311.04254*.
- Duong Tran Dinh, Pham Ngoc Hung, and Tung Nguyen Duy. 2018. A method for automated user interface testing of windows-based applications. In *Proceedings of the 9th International Symposium on Information and Communication Technology*, pages 337–343.
- Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, et al. 2023. CogAgent: A visual language model for gui agents. *arXiv preprint arXiv:2312.08914*.
- Yuxuan Jiang, Chaoyun Zhang, Shilin He, Zhihao Yang, Minghua Ma, Si Qin, Yu Kang, Yingnong Dang, Saravan Rajmohan, Qingwei Lin, et al. 2023. Xpert: Empowering incident management with query recommendations via large language models. *arXiv preprint arXiv:2312.11988*.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. 2023. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026.
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. 2023. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*.
- Bo Qiao, Liquan Li, Xu Zhang, Shilin He, Yu Kang, Chaoyun Zhang, Fangkai Yang, Hang Dong, Jue Zhang, Lu Wang, et al. 2023. TaskWeaver: A code-first agent framework. *arXiv preprint arXiv:2311.17541*.
- Rudolf Ramler, Thomas Wetzlmaier, and Robert Hoschek. 2018. Gui scalability issues of windows desktop applications and how to find them. In *Companion Proceedings for the ISSTA/ECOOOP 2018 Workshops*, pages 63–67.
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Noah Shinn, Beck Labash, and Ashwin Gopinath. 2023. Reflexion: an autonomous agent with dynamic memory and self-reflection. *arXiv preprint arXiv:2303.11366*.
- William Stallings. 2005. The windows operating system. *Operating Systems: Internals and Design Principles*.
- Weihao Tan, Ziluo Ding, Wentao Zhang, Boyu Li, Bohan Zhou, Junpeng Yue, Haochong Xia, Jiechuan Jiang, Longtao Zheng, Xinrun Xu, et al. 2024. Towards general computer control: A multimodal agent for red dead redemption ii as a case study. *arXiv preprint arXiv:2403.03186*.
- Junyang Wang, Haiyang Xu, Jiabo Ye, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang. 2024. Mobile-Agent: Autonomous multi-modal mobile device agent with visual perception. *arXiv preprint arXiv:2401.16158*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Chaoyi Wu, Jiayu Lei, Qiaoyu Zheng, Weike Zhao, Weixiong Lin, Xiaoman Zhang, Xiao Zhou, Ziheng Zhao, Ya Zhang, Yanfeng Wang, et al. 2023. Can GPT-4V (ision) serve medical applications? case studies on GPT-4V for multimodal medical diagnosis. *arXiv preprint arXiv:2310.09909*.
- Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao Yu, and Lingpeng Kong. 2024. Os-copilot: Towards generalist computer agents with self-improvement. *arXiv preprint arXiv:2402.07456*.

- Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. 2023a. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. *arXiv preprint arXiv:2310.11441*.
- Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. 2023b. AppAgent: Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771*.
- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023c. The dawn of Imms: Preliminary explorations with GPT-4V (ision). *arXiv preprint arXiv:2309.17421*, 9(1):1.
- Xinlu Zhang, Yujie Lu, Weizhi Wang, An Yan, Jun Yan, Lianke Qin, Heng Wang, Xifeng Yan, William Yang Wang, and Linda Ruth Petzold. 2023. GPT-4V (ision) as a generalist evaluator for vision-language tasks. *arXiv preprint arXiv:2311.01361*.
- Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. 2024. GPT-4V (ision) is a generalist web agent, if grounded. *arXiv preprint arXiv:2401.01614*.

A More Details on Control Automator

A.1 Example of Control Annotation with SOM

We show an example of control annotation with SOM in Figure 6. Different control types are annotated with different colors. For regions or controls that are not detected by the UI Automation (UIA) framework but still need to be operated on, UFO estimates their coordinates and performs the necessary actions, such as executing a `Click`. This allows UFO to handle scenarios where UIA does not identify certain controls, ensuring comprehensive interaction with the application.

A.2 Supported Action

For the specific functions applied to the control, we have chosen common and widely used mouse operations supported by `pywinauto`, as well as developed customized actions. These actions include `Click`, `SetText`, `GetText`, and `Scroll`, `Annotate` and `Summary`, as detailed below:

- `Click`: Clicking the control item, or a specific coordinate (x, y) on the screen with the mouse, with options for left or right clicks, single or double clicks.
- `SetText`: Inputting text into an editable control, mimicking the keyboard behaviors.
- `SetText`: Inputting text into an editable control, mimicking the keyboard behaviors.
- `Annotate`: Capturing a screenshot of the current application window and annotating the control item on the GUI.
- `Summary`: Summarizing the observation of the current application window based on the clean screenshot.
- `GetText`: Retrieving the textual information of the control.
- `Scroll`: Scrolling the control item vertically or horizontally to make hidden content visible.

A.3 Supported Control Types

UFO focuses on the following 10 constrained control types with high relevance, as determined by our analysis. These include `Button`, `Edit`, `TabItem`, `Document`, `ListItem`, `MenuItem`, `TreeItem`, `ComboBox`,

`Hyperlink`, `ScrollBar`. We show a detailed descriptions of these control types in Table 4 in the Appendix. These set can covers majority of relevant controls in the applications, and it is also extendable per request.

A.4 Sensitive Actions

We show examples sensitive action scenarios that will be confirmed by user before execution in Table 5. UFO automatically determines whether an action is sensitive and requires confirmation. Additionally, users have the option to customize their own sensitive actions that also require confirmation before execution.

B More Details of Experiment Setting

B.1 Parameter Setting

For UFO and all baseline methods, including GPT-3.5, GPT-4, and the agentic baselines, we use GPT-4V as the inference engine. To minimize randomness, we set the parameters `temperature` and `top_n` to 0. For the agentic baselines, we utilized the code from the original repositories with minimal modifications.

B.2 Evaluation Setting

Given the absence of a gold standard or existing open benchmark for evaluating task completion in WindowsBench, human evaluation is employed as the most reliable assessment method. The diversity of tasks, which span multiple applications and types, complicates automated testing and necessitates a case-by-case approach. To ensure a fair and unbiased evaluation, three experienced users of the nine applications included in WindowsBench independently assessed each task. The success and safeguard rates were determined by majority vote, while the completion rate was averaged across the three participants. We obtained explicit consent from all users to utilize their data for the evaluation purposes.

We conducted each test for all frameworks, including UFO and its baselines, three times and selected the best-performing trajectory for evaluation. This approach helps mitigate the randomness inherent in LLM predictions and minimizes the impact of minor variations in the operating environment.

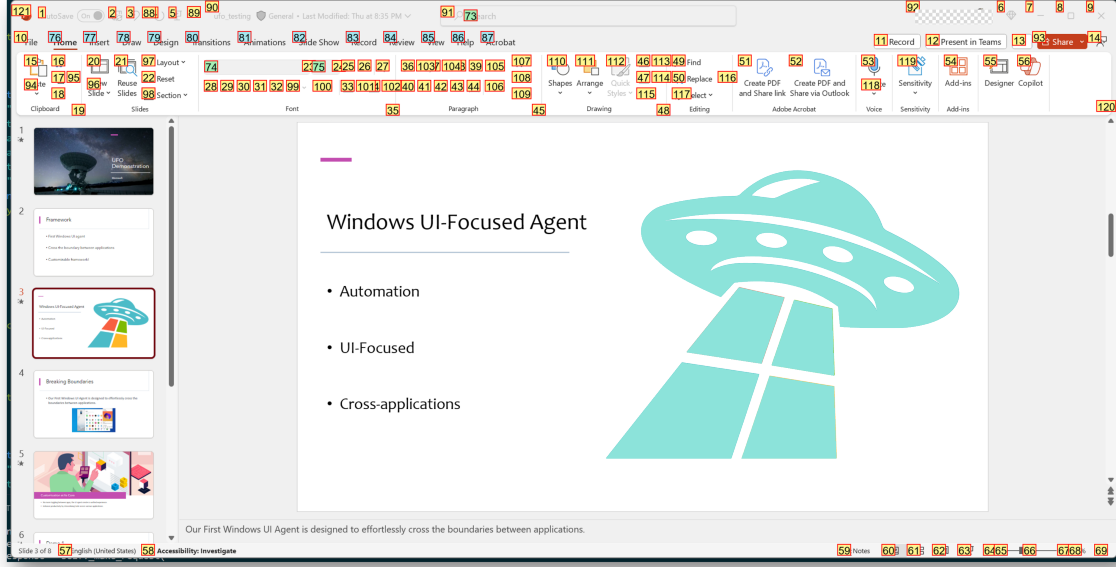


Figure 6: An example of the annotated PowerPoint GUI with information provided by pywinauto. Different colors of annotations represent different control types.

C Details of WindowsBench and Evaluations Breakdown

In Table 6, 7, and 8, we present the complete user requests included in WindowsBench, along with the detailed results achieved by UFO. These requests span across nine different popular Windows applications, incorporating various commonly used functions. Requests with follow-up tasks are sequentially numbered for clarity. In the Safeguard column, “-” denotes that the request is not sensitive and does not require user confirmation. “✓” indicates that the request is sensitive, and UFO successfully activates the safeguard for user confirmation, while “✗” signifies that the safeguard fails to trigger for a sensitive request. In the success column, “✓” denotes that UFO completes the request successfully, while “✗” indicates a failure to fulfill the request.

D Potential Risk

Noted that in the deployment of AI agents capable of executing user requests across applications on Windows OS like UFO, several potential risks warrant careful consideration. These risks can be categorized into security vulnerabilities, operational issues, and ethical concerns:

D.1 Security Vulnerabilities

The use of AI agents to interact with system applications through mouse and keyboard inputs introduces several security risks, since safeguards employed in UFO may not always take effect. Firstly, the agents can be susceptible to exploitation if not properly secured, potentially providing malicious actors with control over the system. For instance, an adversary who gains access to the AI agent could exploit its capabilities to perform unauthorized actions, such as data exfiltration or system modification. Moreover, the agents might inadvertently interact with sensitive information or system settings, increasing the risk of accidental exposure or corruption of critical data.

D.2 Operational Issues

The agents that automate interaction with applications are vulnerable to operational inconsistencies. These agents rely on the assumption that the user interface and application behavior remain constant. However, changes in application updates, variations in user interface design, or unexpected application states can lead to failures or erroneous execution of tasks. Additionally, the reliance on input devices such as the mouse and keyboard may limit the agent’s ability to adapt to diverse or non-standard application environments, thereby reducing its robustness and reliability in handling com-

Table 4: The detailed descriptions of control types supported by UFO.

Control Type	Description
Button	A button is a user interface element that users can interact with to trigger an action. Clicking the button typically initiates a specific operation or command.
Edit	An edit control allows users to input and edit text or numeric data. It is commonly used for fields where users can type information, such as textboxes or search bars.
TabItem	A tab item is part of a tab control, organizing content into multiple pages. Users can switch between different tab items to access distinct sets of information or functionalities.
Document	A document control represents a document or a page in a document-view architecture. It is often used to display and manage documents or large blocks of text.
ListItem	A list item is an element within a list control, presenting data in a list format. Users can select and interact with individual items within the list.
MenuItem	A menu item is part of a menu control, providing a list of commands or options. Users can click on menu items to trigger specific actions or navigate through application features.
TreeItem	A tree item is a node within a tree control, organizing information in a hierarchical structure. Users can expand or collapse tree items to navigate through a hierarchical set of data.
ComboBox	A ComboBox is a combination of a text box and a drop-down list. It allows users to either type a value directly into the text box or choose from a predefined list by opening the drop-down menu.
Hyperlink	A Hyperlink enables users to navigate to another location or resource. They are often used to provide easy access to external websites, documents, or specific sections within an application.
ScrollBar	A scroll bar allows users to scroll through content that is larger than the visible area. It provides a way to navigate vertically or horizontally within a window or control.

plex tasks.

D.3 Ethical Concerns

The deployment of AI agents in user interactions raises ethical considerations related to user autonomy and privacy. Automated systems that perform tasks on behalf of users could potentially infringe upon user agency if not carefully designed to respect user intent. Ethical guidelines must be established to ensure that the AI agent operates transparently, with user consent and oversight, to mitigate risks of misuse or unintended consequences.

Table 5: An incomplete list of sensitive actions considered in UFO.

Sensitive Action	Description
Send Action for a Message or Email	Initiating the “Send” action, such as clicking the send button, is considered sensitive as the dispatched message or email becomes irretrievable.
Deleting or Modifying Files and Folders	Operations involving the deletion or modification of files and folders, especially those situated in critical system directories or containing vital user data.
Closing a Window or Application	Closing a window or application is flagged as sensitive due to the potential for data loss or system crashes.
Accessing Webcam or Microphone	Accessing the webcam or microphone without explicit user consent is identified as sensitive to address privacy concerns.
Installing or Uninstalling Software	Actions related to installing or uninstalling software applications are marked as sensitive due to their impact on system configuration and potential security risks.
Browser History or Password Retrieval	Retrieving sensitive user data such as browser history or stored passwords is identified as a sensitive action, posing potential privacy leaks.

Table 6: The requests in WindowsBench and detailed results achieved by UFO (Part-I).

Request	Application	Safeguard	Step	Success	Completion Rate
I am Zac. Draft an email, to [email address 1] and cc [email address 2], to thanks for his contribution on the VLDB paper on using diffusion models for time series anomaly detection. Don't send it out.	Outlook	-	6	✓	100%
Search 'Spring Festival' in my mail box, and open the second returned email.	Outlook	-	4	✓	100%
Delete the first junk email.	Outlook	✓	4	✓	100%
Please download and save the pdf attachment to the local in the first sent email.	Outlook	-	10	✓	70%
(1) Draft an email, to {email address 1} and to ask him if he can come to the meeting at 3:00 pm. Don't send it out. (2) Add some content to the email body to tell him the meeting is very important, and cc {email address 2} as well. (3) Send the email now.	Outlook	✓	10	✓	100%
Give a detailed description of the components and workflows of the TaskWeaver architecture in the image.	Photos	-	2	✓	100%
Open the image of the LLM training and rotate it for 2 times.	Photos	-	3	✓	100%
Open the image of the LLM training and delete it.	Photos	✓	4	✓	100%
Open the mouse image, then loop over the next two images, without closing them. Summarize them one by one.	Photos	-	5	✗	83.3%
(1) Zoom in the image of the autogen for 2 times to make it clearer, then describe it with detailed workflow and components. (2) It looks too big. Just zoom it to fit the screen.	Photos	-	6	✓	100%
Help me quickly remove all notes in the slide of the ufo_testing, without looping through each slide one-by-one.	PowerPoint	✓	8	✓	100%
Clear the recording on current page for ppt.	PowerPoint	✗	3	✓	100%
Please add Morph transition to the current page and the next page of the ufo_testing.ppt	PowerPoint	-	6	✓	83.3%
Summarize scripts for all pages in the ppt one by one to help with my preseration.	PowerPoint	-	7	✗	85.7%
(1) Apply the first format generated by the Designer to the current ppt.(2) This format is not beautiful. Select a better format and explain why you select it.	PowerPoint	-	8	✓	75%
Check spelling and grammar of the current meeting note.	Word	-	2	✓	100%
Please change the theme of the meeting note into 'Organic'.	Word	-	4	✓	100%
Save meeting note as an Adobe PDF to the local.	Word	-	6	✓	83.3%
Add a cover page to the meeting note. Choose one you think is beautiful.	Word	-	5	✓	100%
(1) Please change to a better-looking format and color to beautify the current page of the meeting note. (2) Do more to make it look even better.	Word	-	10	✓	80%

Table 7: The requests in WindowsBench and detailed results achieved by UFO (Part-II).

Request	Application	Safeguard	Step	Success	Completion Rate
Close all pdf that are currently open.	Adobe Acrobat	✓	3	✓	100%
Cascade the current pdf windows and close the first one.	Adobe Acrobat	✓	3	✗	60%
Visually find and present to me the timeline and milestones in the current meeting note.	Adobe Acrobat	-	4	✗	50%
(1) Read the imdiffuion paper. What is the main contribution of it? (2) Find the first figure of the paper. What does it show?	Adobe Acrobat	-	5	✓	100%
(1) Comprehend the current pdf, what is the difference between phase 1 and phase 2? (2) What is the difference of the Application Selection Agent and the Action Selection Agent? (3) Base on the workflow pdf, is there anything can be improved for the agents? (4) Thanks, this is actually the workflow of yourself. Do you think there are other things can be improved for you? I will implement your thought on it. (5) Great. Do you think there are other things that can be improved visually for the figure itself in the pdf, to provide better illustration for other readers?	Adobe Acrobat	-	6	✓	83.3%
Navigate to /Desktop/mix/screen and open s0.png by only double clicking.	File Explorer	-	6	✓	83.3%
Delete all files individually in the current figure folder.	File Explorer	✓	6	✓	100%
Copy the AppAgent file in the folder to the Document folder.	File Explorer	-	5	✓	60%
Create a new txt file in the current folder.	File Explorer	-	3	✓	100%
(1) Open all images that are related to dragon in the test folder. (2) Summarize in detailed what you have seen in the image d4.	File Explorer	-	4	✓	100%
Tell me the changelog of the {code base name} repo in their readme base on your observation.	Visual Studio Code	-	5	✓	80%
Find the use of 'print_with_color' in the script folder in the {code base name} repo.	Visual Studio Code	-	4	✓	100%
Download the Docker extension in the {code base name} repo.	Visual Studio Code	-	4	✓	100%
Create a new terminal at the {code base name} repo.	Visual Studio Code	-	2	✓	100%
(1) Read through carefully and review the code in the utils.py in the {code base name} and identify any potential bugs or aspects can be improved. (2) What about in the draw_bbox_multi function?	Visual Studio Code	-	5	✗	40%
Please send a 'smile' emoticon to the current chatbox at Wechat.	WeChat	✗	4	✓	100%
Delete the current chatbox of 'File Transfer' on Wechat.	WeChat	✓	3	✓	100%
Please like {user name}'s post at Moment.	WeChat	-	3	✓	100%
Bring the first picture to the front in my "favourite" at WeChat and decribe it.	WeChat	-	5	✓	100%
(1) Open the {account name} Official Account at WeChat. (2) Send a message to 'File Transfer' to introduce this account.	WeChat	✓	10	✓	90%

Table 8: The requests in WindowsBench and detailed results achieved by UFO (Part-III).

Request	Application	Safeguard	Step	Success	Completion Rate
How many total citations does Geoffrey Hinton have currently?	Edge Browser	-	4	✓	100%
Post "It's a good day." on my Twitter.	Edge Browser	✓	6	✓	100%
Download the Imdiffusion repo as zip.	Edge Browser	-	7	✓	100%
Change the theme of the browser to the icy mint.	Edge Browser	-	5	✓	100%
(1) Find and navigate to {name}'s homepage at Microsoft. (2) Print this page in color.	Edge Browser	-	5	✗	60%
My name is Zac. Please read the meeting note in meeting_note.docx to identify all action items, and include the detailed description LLM training workflow in the LLM-training.png, to compose an new email of these contents. Sent the full email to our leader Hidan via {email address} to ask for his review.	Word, Photos, Outlook	-	10	✓	100%
Search for and read through the latest news about Microsoft and summarize it, then send it to 'File Transfer' on WeChat.	Edge Browser, WeChat	✓	9	✓	100%
Open the image of ufo_rv at UFO-windows/logo/, summarize its content and use the summary to search for a similar image on Google.	File Explorer, Edge Browser	-	10	✓	90%
Observe the overview.pdf, then download the most similar paper base on this architecture from the Internet base on your understanding, and sent it to the 'File Transfer' on WeChat.	Adobe Acrobat, Edge Browser, WeChat	-	10	✗	25%
(1) Read the paper.pptx, summarize the paper title presented on the slide, then search and open the paper on the Internet. (2) Can you summarize this paper and download its PDF version for me?	PowerPoint, Edge Browser	-	10	✓	100%

E Performance Breakdown of GPT-3.5 and GPT-4

In Tables 9 and 10, we present the detailed performance breakdown for GPT-3.5 and GPT-4, respectively. It is evident that GPT-4 significantly outperforms GPT-3.5, aligning with our expectations. Notably, GPT-4 exhibits superior performance, particularly in tasks related to `Visual Studio Code`, where GPT-3.5 frequently struggles to select the correct application initially, leading to its overall lower performance. Another noteworthy instance is seen in `Edge Browser`, where GPT-3.5 often overlooks crucial steps, resulting in request failures. In general, both baselines exhibit inconsistent performance across different applications, and their overall efficacy falls significantly short when compared to UFO.

F Additional Case Study

In this section, we present six additional case studies to illustrate the efficacy of UFO in performing diverse tasks across the Windows operating system. These case studies encompass both single-application scenarios and instances where UFO seamlessly transitions between multiple applications.

F.1 Deleting All Notes on a PowerPoint Presentation

In Figure 7, we tasked UFO with the request: “Help me quickly remove all notes in the slide of the `ufo_testing`.”. This request is a common scenario when users want a clean version of a slide without any accompanying notes. Traditional methods involve manually deleting notes page by page, a tedious and time-consuming process for lengthy slides.

However, UFO identifies a shortcut, simplifying the entire procedure. Analyzing its initial plan, UFO proposes the use of the ‘Remove All Presentation Notes’ function, which is a feature often overlooked by PowerPoint users due to its hidden placement. The automation process commences with UFO navigating to the “File” tab, providing access to the backstage view. Subsequently, it smoothly transitions to the “Info” menu, then clicks the ‘Check for Issues’ button and selects ‘Inspect Document’. Once inspected, the hidden gem, ‘Remove All Presentation Notes’ is revealed at the bottom of the menu. UFO scrolls down, locates the function, and initiates the click. Given the sen-

sitive nature of deleting notes, UFO implements its safeguard feature, seeking user confirmation. Upon confirmation, all notes vanish with a single click. This demonstration highlights how UFO empowers users to work smarter, not harder, establishing itself as an intelligent assistant.

F.2 Reading a PDF

UFO also demonstrates its effectiveness in reading and comprehending a PDF file. In Figure 8, we request UFO to “(1) Read the imdiffusion paper. What is the main contribution of it? (2) Find the first figure of the paper. What does it show?”. The imdiffusion paper (Chen et al., 2023) delves into complex technology related to anomaly detection and diffusion models. Remarkably, UFO showcases its understanding of the paper by summarizing its abstract using OCR in vision input, accurately capturing its main contribution. In a follow-up request, we ask UFO to find and explain the first figure. UFO navigates through the PDF, locates the initial figure, and delivers a detailed explanation, comparing the forecasting, reconstruction, and imputation approaches in anomaly detection. This showcases the remarkable OCR ability of UFO in understanding both text and figures. Below shows UFO’s summary on the paper and the description of the first figure:

Paper Summary: The ‘ImDiffusion’ paper’s abstract discusses the importance of anomaly detection in multivariate time series data and the challenges associated with existing approaches. The paper introduces ‘ImDiffusion’, a novel approach that combines time series imputation and diffusion models to achieve accurate and robust anomaly detection. The method is designed to handle complex multivariate time series data, enabling the imputation of missing and irregular values in the time series. This approach reportedly improves the accuracy of anomaly detection by leveraging diffusion models as a preprocessing step to accurately reconstruct the anomaly detection time series.

Table 9: The detailed performance breakdown achieved by GPT-3.5 (Human Surrogate) on WindowsBench.

Application	Success	Step	Completion Rate	Safeguard Rate
Outlook	60.0%	7.6	76.2%	100.0%
Photos	60.0%	5.6	35.7%	100.0%
PowerPoint	20.0%	11.2	20.6%	0.0%
Word	20.0%	7.6	25.0%	-
Adobe Acrobat	0.0%	7.4	5.3%	50.0%
File Explorer	40.0%	7.0	53.3%	0.0%
Visual Studio Code	0.0%	7.6	0.0%	-
WeChat	40.0%	6.0	48.4%	33.3%
Edge Browser	0.0%	7.6	32.0%	100.0%
Cross-Application	0.0%	11.0	18.8%	50.0%

Table 10: The detailed performance breakdown achieved by GPT-4 (Human Surrogate) on WindowsBench.

Application	Success	Step	Completion Rate	Safeguard Rate
Outlook	100.0%	8.4	73.9%	0.0%
Photos	40.0%	7.0	32.7%	100.0%
PowerPoint	40.0%	10.4	35.2%	50.0%
Word	20.0%	9.2	15.3%	-
Adobe Acrobat	0.0%	7.6	40.2%	50.0%
File Explorer	80.0%	6.2	63.4%	100.0%
Visual Studio Code	40.0%	7.4	40.3%	-
WeChat	40.0%	6.2	68.0%	66.7%
Edge Browser	60.0%	8.2	58.8%	100.0%
Cross-Application	0.0%	13.8	49.7%	50.0%

Figure Description (Part-I): Figure 1: Examples of reconstruction, forecasting and imputation modeling of time series for anomaly detection. To address the challenges and overcome the limitations of existing approaches, we propose a novel anomaly detection method - ImDiffusion. This detector combines the use of time series imputation and diffusion models to achieve accurate and robust anomaly detection. ImDiffusion employs dedicated gating data that makes the diffusion models extra cautiously model the time points, helping utilities see more detail, to recognize valued and serious issues of anomalies through the challenges of detection. The paper introduces a novel framework based on the combination of time series imputation and diffusion models that effectively leverages neighboring values in the time series as additional conditional information,

Figure Description (Part-II): enabling a more accurate modeling of the temporal and inter-related dependencies present in multi-variate data. Secondly, the reference performance threshold of the calibration test of the counter prediction from then exchanging values helps to focus the detection in outliers, Figure 1 presents an example in which forecasting, reconstruction, and imputation models are employed to predict a time series using diffusion models. Observe that the imputation-based methods exhibit the highest prediction performance especially on normal periods, thus only the imputation method successfully identifies the period of anomaly. We therefore employ time series imputation for accurate self-supervised modeling of time series, which forms the foundation of our proposed ImDiffusion framework.

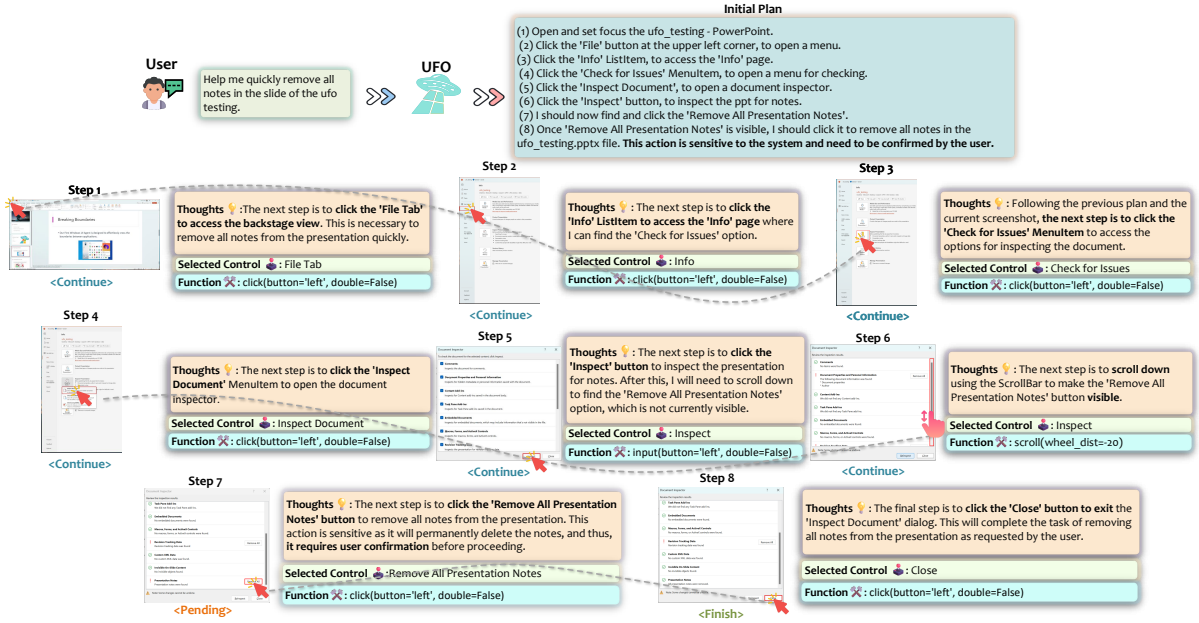


Figure 7: An example of UFO removing all PPT notes.

F.3 Designing a PowerPoint Slide

In Figure 9, we present another example where we make two requests to UFO: “(1) Apply the first format generated by the Designer to the current ppt. (2) This format is not beautiful. Select a better format and explain why you selected it.” The first step is relatively straightforward, where UFO successfully navigates to the ‘Designer’ Pane and applies the desired format to the slide. The second request, however, depends on UFO’s aesthetic judgment. UFO performs this task effectively and provides a reasonable explanation: “I have selected the design suggestion because it provides a clear visual hierarchy by placing the descriptive picture at the top, drawing immediate attention. The title and text are positioned at the bottom, making it easy for the audience to follow the flow of information. The contrasting background colors and font colors enhance readability and focus on the content”. This showcases the remarkable visual ability of UFO, as it can complete open-ended requests based on its own aesthetic sense.

F.4 Downloading an Extension for VSCode

In Figure 10, we present an example where UFO effectively fulfills a user request in Visual Studio Code: “Download the Docker extension in the repo.”. Although this is a relatively simple task, we aimed to test UFO’s performance in operating on a dark mode GUI. Remarkably, UFO completes this request effortlessly by clicking the Extension but-

ton, inputting “Docker” into the search box, and accurately clicking the “Install” button for the Docker Extension. With just one click, the task is successfully completed. This showcases UFO’s capability to operate on less popular applications, even when they are in dark mode.

F.5 Post a Twitter

In Figure 11, we shift our focus to the Edge Web browser to assess UFO’s capability in operating over one of the most popular application types on Windows OS. Specifically, we request UFO to perform the task of “Post ‘It’s a good day.’ on my Twitter.”, a common yet intricate request for an agent. The observation reveals UFO’s seamless execution, where it inputs the Twitter address into the browser’s address bar, navigates to the link perceived as the main Twitter page, identifies the Post button on the Twitter page, and clicks it. This decision, although not necessarily straightforward, allows UFO to input the required text for the tweet. The safeguard activates before sending, prompting user confirmation. The successful posting of the tweet demonstrates UFO’s sophisticated ability to operate on web browsers, a category of applications widely used on Windows OS.

F.6 Sending the News

In a cross-application example, we illustrate how UFO can gather news from the Internet and share it on WeChat in Figure 12. Specifically, we issue the command: “Search for and read through the latest

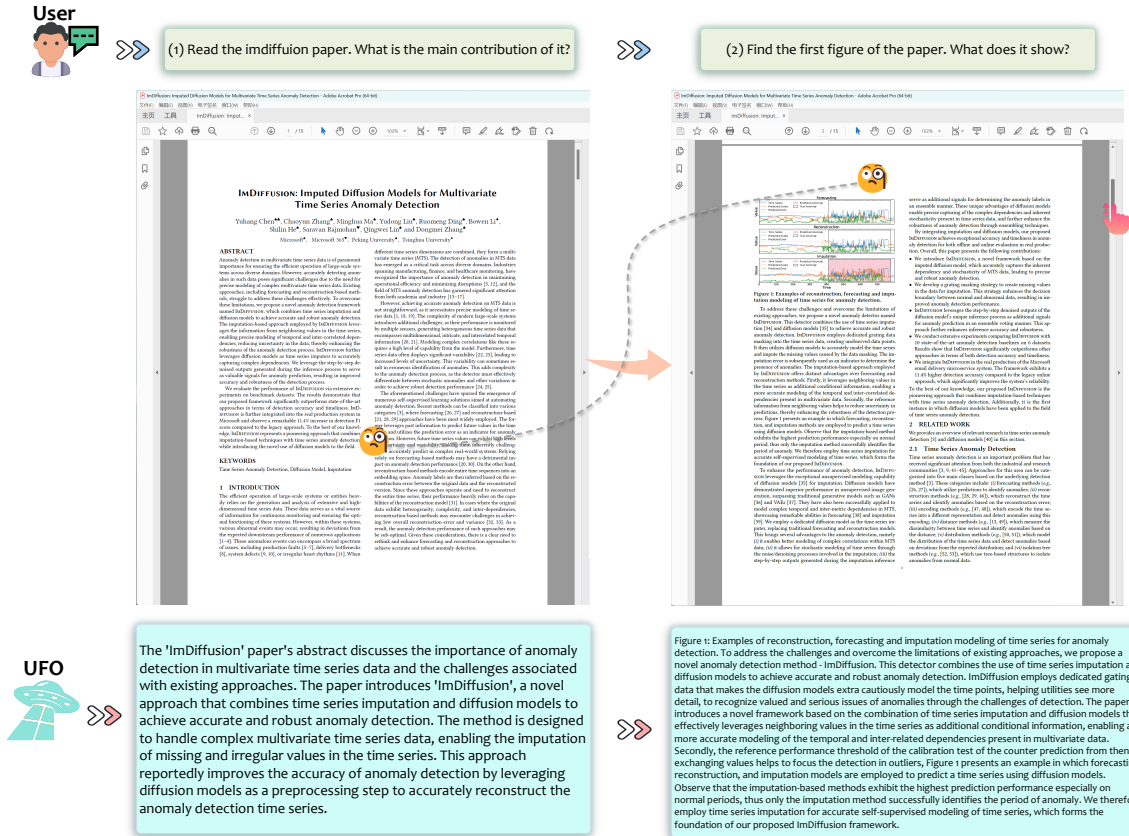


Figure 8: UFO completes user request: “(1) Read the imdiffusion paper. What is the main contribution of it? (2) Find the first figure of the paper. What does it show?”

news about Microsoft and summarize it, then send it to 'File Transfer' on WeChat". UFO adeptly inputs the query "latest news about Microsoft" into the Google search bar on Edge Browser, initiates the search, and opens the first link in the search results. The page contains multiple news panes, and UFO skillfully summarizes the content using its visual OCR ability, logging the information into its memory. Subsequently, UFO opens WeChat, locates the designated chatbox, and inputs the summarized news retrieved from its memory. Upon user confirmation triggered by the safeguard, UFO sends the news. We show the news sent below.

The Microsoft News page displays several articles. The headlines include: 'Excited to announce the general availability of Copilot for Sales and Copilot for Service, as we continue to extend Copilot to every role and function.', 'Microsoft announces quarterly earnings release date', 'With Copilot Pro, we're helping even more people supercharge their creativity and productivity by unlocking Copilot in Microsoft 365 apps, providing access to the very latest models 2014 and more.', and 'Microsoft unveils new generative AI and data solutions across the shopper journey, offering copilot experiences through Microsoft Cloud for Retail'.

This example once again underscores UFO's ability to seamlessly transition between different applications, allowing it to effectively and safely complete long-term, complex tasks. Such demonstrations position UFO as an advanced and compelling agent for Windows OS.

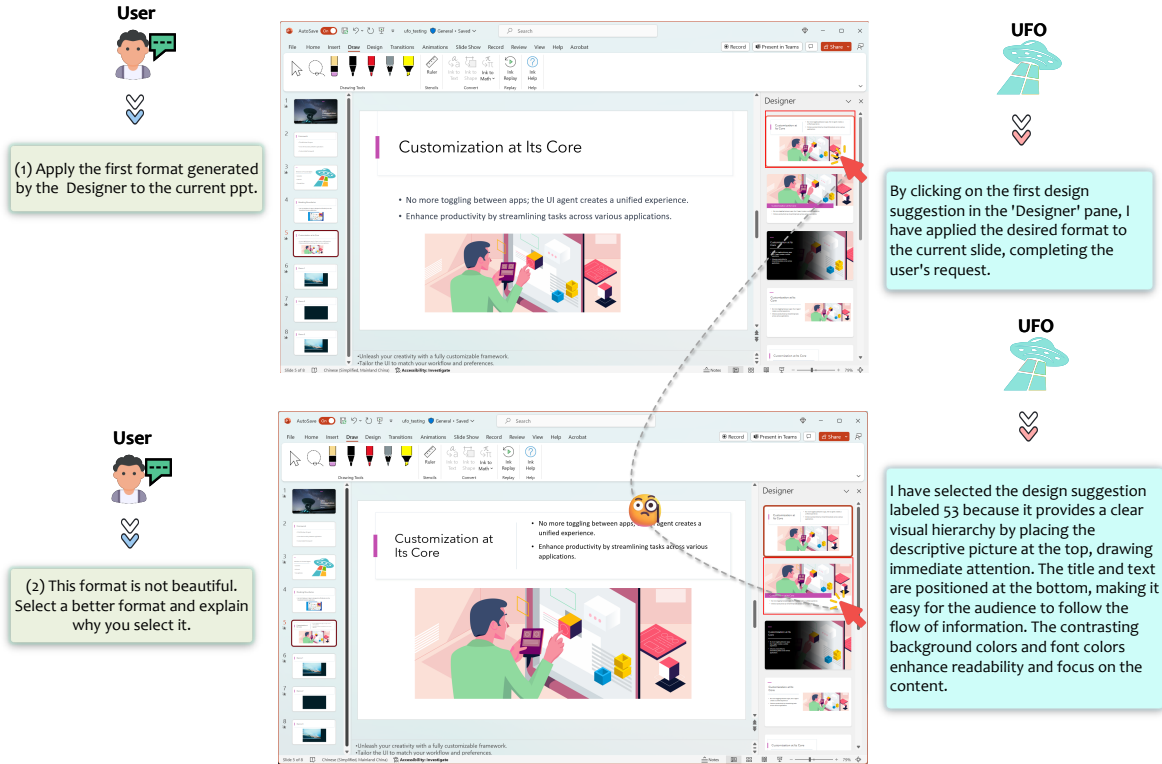


Figure 9: UFO completes user request: “(1) Apply the first format generated by the Designer to the current ppt. (2) This format is not beautiful. Select a better format and explain why you select it.”

F.7 Search Paper in a PowerPoint Slide and Summarize

Lastly, in Figure 13, we present an intricate cross-application example with multiple sub-requests where UFO is tasked with finding a paper presented in a slide on the Internet, summarizing the paper, and downloading it. The user’s request is as follows: “(1) Read the paper.pptx, summarize the paper title presented on the slide, then search and open the paper on the Internet. (2) Can you summarize this paper and download its PDF version for me?” This request involves multiple instances of OCR usage by UFO.

To fulfill the first sub-request, UFO carefully observes the title page of the slide, successfully identifies the paper titled “Xpert: Empowering Incident Management with Query Recommendations via Large Language Models.” Subsequently, it opens a Google page in the Edge Browser, inputs the title into the search bar, and locates the paper as the second result by skipping the first ad link. UFO clicks and opens the link on arXiv, completing the first request.

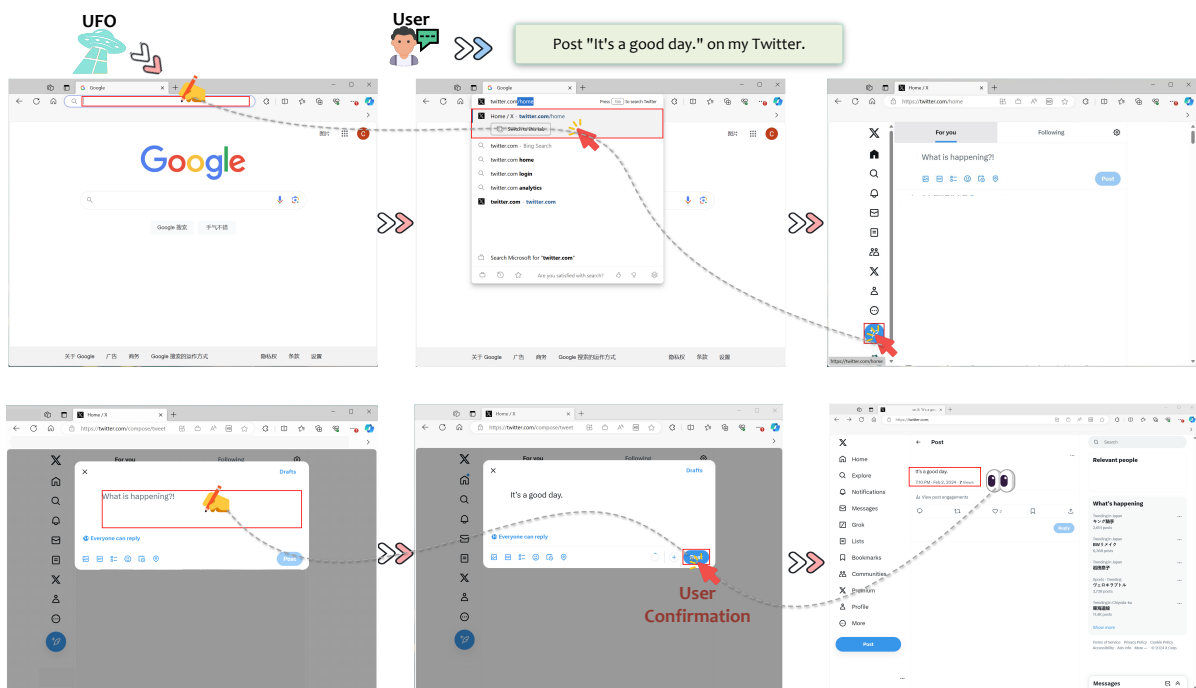
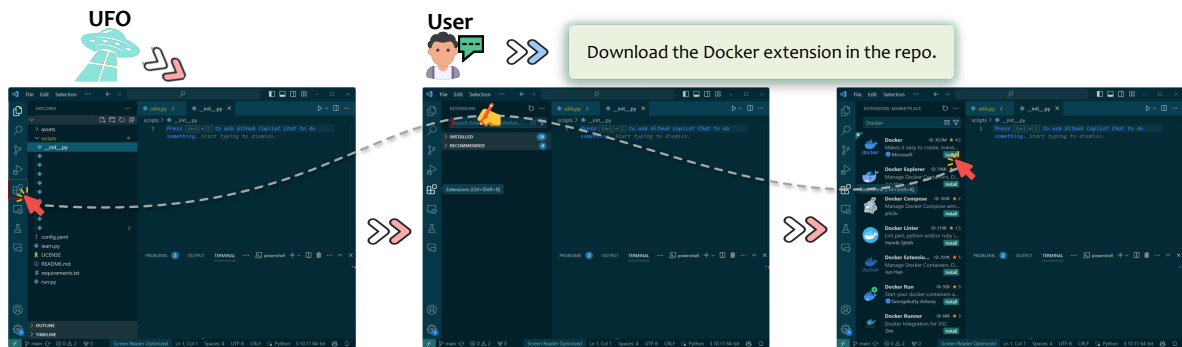
Upon receiving the second sub-request to summarize the paper and download the PDF version, UFO clicks the “Download as PDF” option and

uses its OCR ability to generate the following summary:

Part I: The paper titled ‘Xpert: Empowering Incident Management with Query Recommendations via Large Language Models’ discusses the role of large-scale cloud systems in modern IT infrastructure and the impact of incidents occurring within these systems. The paper introduces Xpert, a machine learning framework that automates Kusto Query Language (KQL) recommendations, tailored to new incidents. By leveraging historical incident data and large language models, Xpert generates customized KQL queries, enabling a thorough evaluation of query quality from three comprehensive perspectives.

Part II: The paper details the deployment of Xpert in the real production environment of a large-scale incident management system in Microsoft, demonstrating its effectiveness in supporting incident management.

This summary aligns perfectly with the paper’s con-



tribution as outlined in (Jiang et al., 2023), showcasing UFO’s excellent ability to operate across different applications and its remarkable OCR capabilities.

G The Use of AI Assistance

We use ChatGPT to polish the content of the paper.

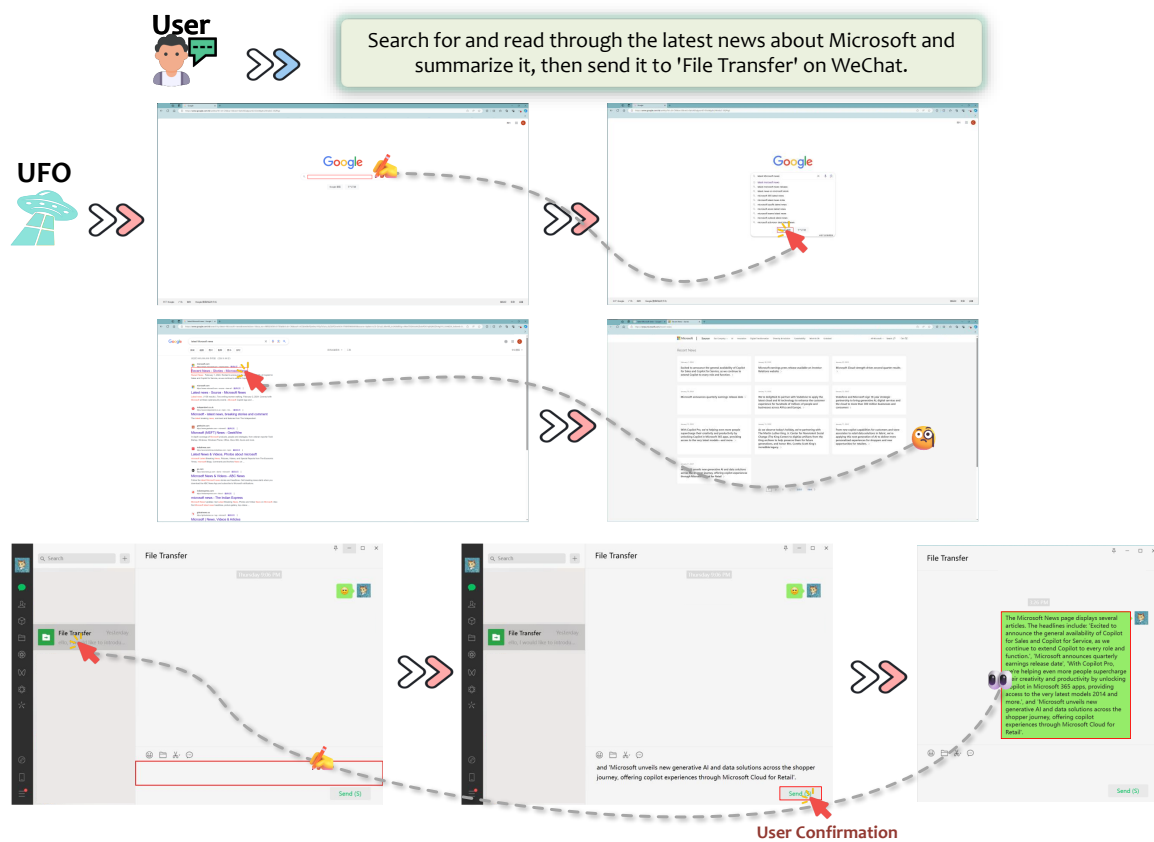


Figure 12: UFO completes user request: “Search for and read through the latest news about Microsoft and summarize it, then send it to 'File Transfer' on WeChat.”.

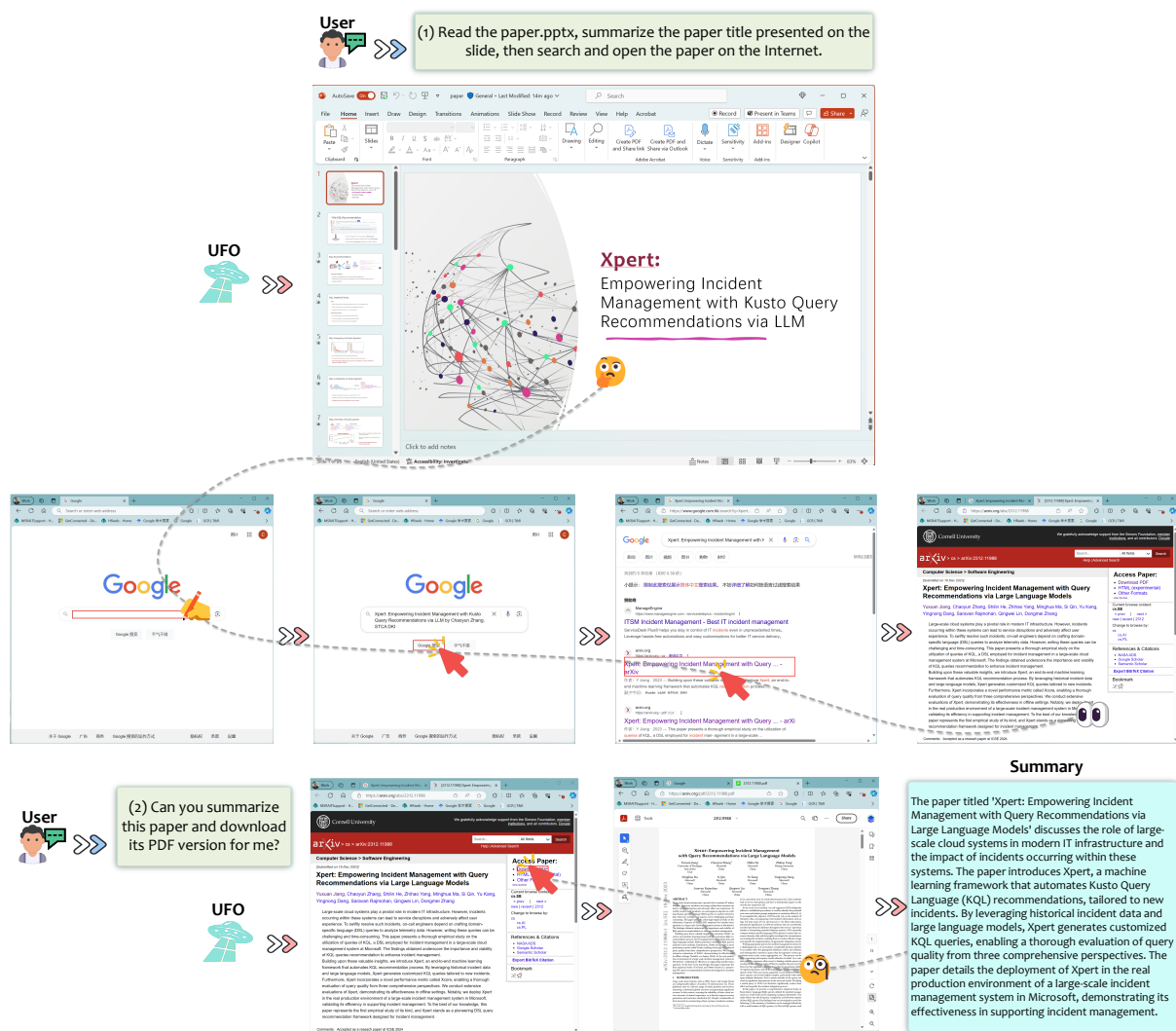


Figure 13: UFO completes user request: “(1) Read the paper.pptx, summarize the paper title presented on the slide, then search and open the paper on the Internet. (2) Can you summarize this paper and download its PDF version for me?”.