

Learning from Online Videos at Inference Time for Computer-Use Agents

Anonymous authors

Paper under double-blind review

Abstract

Computer-use agents can operate computers and automate laborious tasks, but despite recent rapid progress, they still lag behind human users, especially when tasks require domain-specific procedural knowledge about particular applications, platforms, and multi-step workflows. Humans can bridge this gap by watching video tutorials: we search, skim, and selectively imitate short segments that match our current subgoal. In this paper, we study how to enable computer-use agents to learn from online videos at inference time effectively. We propose a framework that retrieves and filters tutorial videos, converts them into structured demonstration trajectories, and dynamically selects trajectories as in-context guidance during execution. Particularly, using a VLM, we infer UI actions, segment videos into short subsequences of actions, and assign each subsequence a textual objective. At inference time, a two-stage selection mechanism dynamically chooses a single trajectory to add in context at each step, focusing the agent on the most helpful local guidance for its next decision. Experiments on two widely used benchmarks show that our framework consistently outperforms strong base agents and variants that use only textual tutorials or transcripts. Analyses highlight the importance of trajectory segmentation and selection, action filtering, and visual information, suggesting that abundant online videos can be systematically distilled into actionable guidance that improves computer-use agents at inference time.

1 Introduction

Computer-use agents promise to automate labor-intensive digital workflows, reshape how people interact with computers, and boost productivity across everyday and expert tasks (Deng et al., 2023; Xie et al., 2024; Zheng et al., 2024; PBC, 2024; OpenAI, 2025). Recent works have made rapid progress on a wide range of desktop and web tasks, yet a substantial gap remains between their performance and that of human users (Wang et al., 2025a; Agashe et al., 2025b; Wang et al., 2025b). We argue that a key reason for this gap is the agent’s limited access to domain-specific procedural knowledge: how to operate particular applications, cope with idiosyncratic UI conventions, and execute multi-step workflows that depend on software version and platform. While agents continue to improve, progress on these capabilities is constrained by the scarcity of such data in large language models’ (LLMs) training corpus.

On the contrary, humans are usually not limited by the lack of such domain knowledge. When we encounter unfamiliar applications or platforms, we rarely rely on our general knowledge alone. Instead, we consult online video tutorials and quickly learn useful information from the demonstrations. More specifically, we are able to search, skim, and selectively imitate video segments that match our immediate subgoals. These videos are abundant and richly visual, showcasing concrete UI manipulations that text-only resources frequently omit. This observation motivates our central research question: *How can we enable computer-use agents to learn effectively from online videos at inference time?*

Answering this question is challenging for two reasons. First, there is a gap between watching videos and completing tasks as computer-use agents. Videos are continuous streams of frames, whereas agents issue actions as individual events, and they only observe sparse screenshots after completion of each action. Moreover, agents need to generate explicit UI actions, but actions performed in the videos are implicit and

need to be inferred from screen transitions. Bridging from pixel changes in a video to actionable steps is therefore nontrivial. Second, raw videos are long and often contain digressions. During execution, an agent usually needs only a short, locally relevant snippet that matches its current observation and subgoal, so injecting entire videos can distract more than help.

To address these challenges, we introduce an inference-time learning-from-video framework that turns raw video tutorials into compact, actionable demonstrations the agent can consult step by step. Our approach has three components. First, given the task description, the agent generates search queries, searches online videos, and filters them to a small set of relevant tutorials. Second, a processing pipeline converts each filtered video into structured demonstration trajectories. Specifically, an off-the-shelf vision language model (VLM) infers underlying UI actions and synthesizes a concise objective for each subsequence of actions. With additional filtering of low-quality subsequences, this step yields a set of demonstration trajectories, each with a textual objective and a sequence of actions and observations. Third, a dynamic selection method ensures the agent consumes only what matters at the moment of decision. Particularly, before issuing the next action, the agent performs a coarse ranking of demonstrations followed by a detailed inspection of action lists to select a single trajectory that is most helpful for the next action. The selected trajectory is then provided in context to guide the agent’s next decision. This design operationalizes how humans learn from videos: we search, skim, and imitate only the few seconds that matter. By converting videos into short, validated, visually grounded trajectories, the agent receives exactly the local guidance it needs, in the format it can act on.

We evaluate our method on OSWORLD-VERIFIED (Xie et al., 2024) and WEBARENA (Zhou et al., 2024), two widely used benchmarks for desktop and web tasks, respectively. Across both settings, our approach consistently outperforms strong agent baselines that lack video access, as well as variants that use only textual tutorials or video transcripts. For example, on OSWORLD-VERIFIED, our method outperforms the state-of-the-art Jedi framework (Xie et al., 2025b) by 2.1. On WEBARENA, our method outperforms AgentOccam (Yang et al., 2025b) by 4.2. Additional analyses show that our method benefits from more relevant videos, indicating the potential for further performance scale-up. Ablation studies also demonstrate the effectiveness of our designs, such as segmenting videos into short trajectories and dynamically selecting them, as well as filtering out irrelevant actions during video processing. Together, these results suggest that abundant online video tutorials can be systematically and automatically harvested and distilled into compact, actionable guidance that improves the ability of current computer-use agents.

2 Related Works

Computer-use agents. Recent advances in computer-use agents aim to endow LLM- or VLM-based systems with the ability to interact with GUIs on desktops and the web to complete tasks. Early work has focused on constructing realistic benchmarks and environments to support systematic evaluation (Xie et al., 2024; Deng et al., 2023; Zhou et al., 2024; Lin et al., 2024; Rawles et al., 2025; Bonatti et al., 2024; Zhao et al., 2025; Chen et al., 2025). Building on these benchmarks, subsequent systems explore different ways to couple multimodal perception with planning and control, including vision-enabled generalist agents, compositional generalist-specialist architectures, and reinforcement learning pipelines that improve agents’ ability through training (Agashe et al., 2025b; Wang et al., 2025a; Qin et al., 2025; Wu et al., 2024; Yang et al., 2025c; Sun et al., 2025; Tan et al., 2025; Xu et al., 2025b; Xie et al., 2025a; Jia et al., 2025; Gao et al., 2024).

Similar to our work, some recent papers have investigated collecting demonstration trajectories and applying them in context or leveraging them for training. For example, Su et al. (2025) collects and filters trajectories from the agent’s self-exploration in the environment. Ou et al. (2024) extracts trajectories from online tutorials, and Xu et al. (2025a) further simulates execution traces by following retrieved tutorials. However, self-exploration is limited by the agent’s own capability, making it difficult to synthesize complex trajectories and tasks. While online tutorials provide abundant knowledge, most existing work focuses on textual tutorials, leaving the rich visual information in videos under-exploited.

Imitation learning from observations. Our approach is conceptually related to imitation learning from observations in reinforcement learning, where agents learn from state-only demonstrations without access to expert actions (Liu et al., 2018; Aytar et al., 2018; Yu et al., 2018; Torabi et al., 2019; Guo et al., 2019;

Schmeckpeper et al., 2021). In particular, Torabi et al. (2018) has shown that it is possible to train an inverse dynamics model that infers latent actions from state transitions. The trained model is then used to label trajectories in expert demonstrations, which are in turn used for supervised training.

Several recent works have adopted a similar perspective for computer-use agents: tutorial videos provide rich observational demonstrations of how to operate software, but must be converted into actionable steps that agents can execute (Zhang et al., 2025; Jang et al., 2025b; Song et al., 2025; Lu et al., 2025). Specifically, Zhang et al. (2025) crawls online multimodal GUI tutorials (both videos and articles) and converts them into agent trajectories via VLM-based tutorial processing, trajectory generation, and data filtering. The collected trajectories are then used to improve agents by continuous pre-training. Similarly, two concurrent works (Song et al., 2025; Lu et al., 2025) also conduct continuous pre-training from unlabeled videos. To more accurately extract the underlying actions in the video, they further train an inverse dynamics model specialized for the GUI operations.

These works are complementary to ours. They focus on how to accurately label actions and scale up video-derived datasets for offline training, whereas we focus on an inference-time method that does not require additional parameter updates. Given a set of labeled trajectories, our framework studies how to search, filter, and contextualize them so that a computer-use agent can effectively exploit this procedural knowledge on the fly. In addition, our method can potentially be further enhanced if equipped with the specialized inverse dynamics model trained from their works.

3 Methodology

3.1 Pipeline Overview

In this paper, we aim to develop computer-use agents that automatically complete tasks on the desktop. Given an input task, *e.g.*, “*Fill all the blank cells in B1:E30 with the value in the cell above it*”, the initial state, *e.g.*, an Ubuntu desktop with specific files opened, the agent operates on the computer by issuing mouse and keyboard actions. At each step, the agent observes the screenshot of the current desktop and optionally the accessibility tree, and issues the next action by specifying the action type (*e.g.*, click or type) and content (*e.g.*, where to click or what to type). Our goal is to improve the ability of the agent at *inference time* by retrieving relevant online demonstration videos and processing them in a way that allows the agent to learn useful information in context.

To achieve such a goal, we design a pipeline that consists of the following three steps, as illustrated in Figure 1.

1. **Video retrieval:** The agent retrieves and filters relevant online videos that contain useful demonstrations applicable to the testing task.
2. **Video processing:** A VLM extracts underlying actions in the video and summarizes important video segments into demonstration trajectories that contain a description and a sequence of actions and observations.
3. **Video application:** During inference, we design a two-stage selection method to select the most relevant demonstration trajectory at each step. The selected trajectory is provided in context to help the agent better decide the next action.

In the following section, we will describe the details of each step. All prompts used in our pipeline are presented in Appendix A.

3.2 Video Retrieval

Given a computer-use task described in natural language, our first step is to retrieve online video tutorials that plausibly demonstrate procedures relevant to the task. To achieve this, the agent generates queries for each task, searches for videos using those queries, and filters retrieved videos.

Retrieving online videos. We first prompt LLMs to convert the input task into concrete search queries that reflect the application, objects, and intended operation. For the above example, we issue the query:

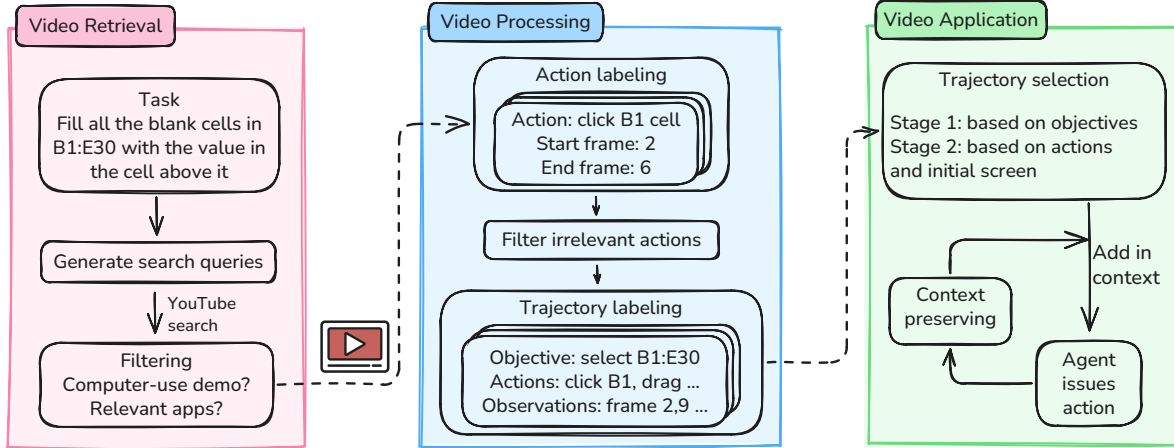


Figure 1: Overview of our inference pipeline. Step 1 outputs a set of videos relevant to the query task. Step 2 yields a set of demonstration trajectories, each with a textual objective and a sequence of actions and observations. Step 3 dynamically selects a single trajectory to add in context before issuing each action, focusing the agent on the most helpful local guidance for its decision.

‘*LibreOffice Calc, fill blank cells with value from above*’. If the task involves multiple applications, we allow the model to generate one query for each application. We then submit each query to the YouTube API¹ and collect the top-50 results per query. Among the retrieved videos, we only keep those that are shorter than 10 minutes and in English.

Filtering retrieved videos. The retrieved videos might be noisy, *e.g.*, a recording of a person speaking without computer operations, or a tutorial on Windows while the task is for Ubuntu, so we filter the retrieved videos to a compact, high-precision set based on LLM assessments.

Specifically, the filtering is done in two steps. ❶ *Coarse selection*. Given the task description, as well as the title and description of all videos, the LLM selects the top-10 videos that are relevant to the task. ❷ *Content verification*. For each kept video, we use a VLM to examine its transcript and a uniform sample of 10 frames spanning the video. The VLM decides whether the video actually shows the application UI and computer operations pertinent to the task. We only keep videos that are deemed helpful.

This stage returns a small set of videos that the VLM has judged as both topically and visually relevant to the task. Multiple videos can be kept for a single task to account for UI and version diversity. A representative result from this stage is a LibreOffice Calc tutorial.²

3.3 Video Processing

Although the filtered videos contain relevant demonstrations for the task, they are not presented in a way that is ready for use by the agent. Particularly, the video contains continuous frames, but the agent only has sparse observations (*i.e.*, screenshots before and after each action). Moreover, the actions performed in the video are not explicitly labeled and must be inferred from screen transitions. To bridge this gap between videos and the agent, we convert each video into structured demonstration trajectories, where each trajectory includes an objective and a sequence of observations and actions.

Action labeling with a VLM. The first step of video processing is to label the underlying actions in the video. To make the labeling efficient, we first identify the frames with screen changes and only annotate actions across these frames. Specifically, we downsample each video with 2 frames per second and locate

¹<https://developers.google.com/youtube/v3>.

²<https://www.youtube.com/watch?v=AeVBHOC1CoA>

changing frames whose pixel value difference from the previous frame is above a certain threshold. We then prompt a VLM to label the actions in the video. To balance the performance and efficiency, we provide the VLM with a context containing 20 frames with 3 frames of overlap between consecutive context clips. The VLM then extracts ① the UI action type (click, type, right-click, drag, *etc.*), ② the target element (*e.g.*, ‘cell B2’, ‘Find All’ button), and ③ start and end frames of the action. Finally, we filter actions based on action types and prompt an LLM to merge the same underlying actions in adjacent clips.

Filtering important and relevant actions. The above step is very likely to output actions that are irrelevant to the task, *e.g.*, cursor movement or hovering without any purpose. To improve the quality of labeled actions, we further prompt an LLM to analyze each action and only keep actions that are important and relevant to the main task in the video based on the title, description, and transcript of the video.

Demonstration trajectory construction. Even after filtering, a video could contain many actions, which is inefficient for the agent to process. In addition, to determine the next action, the agent usually only needs to focus on a local segment of the video instead of looking at the whole video. To allow the agent to focus on a local segment, we further split the video into subsequences of actions and construct a trajectory for each valid subsequence.

Specifically, we take inspiration from the hindsight relabeling method in RL (Andrychowicz et al., 2017; Su et al., 2025) and re-generate a new objective for a trajectory that does not complete the original task. In particular, given a subsequence of observations and actions, we prompt a VLM to generate a natural language objective that this particular subsequence can accomplish. For example, in our above tutorial on LibreOffice Calc, an objective generated for a segment is ‘*Search for all occurrences in the spreadsheet using the ‘Find All’ function in LibreOffice Calc.*’ We repeat such a process for all subsequences of actions in the video, and we only keep the results where the VLM deems the subsequence as complete and there exists a reasonable objective. The outcome of this step is a set of *demonstration trajectories*,

represented as $\left\{ \left(\text{objective}, \{(o_t, a_t)\}_{t=1}^{l_i} \right) \right\}_{i=1}^K$, where o_t is the screenshot at step t , a_t is the corresponding action, l_i is the length of i -th trajectory, and K is the total number of trajectories in the video.

The VLM might generate sub-optimal objectives for the trajectory, where the action sequence does not fully complete the objective. To ensure high quality of constructed trajectories, we further prompt a VLM to filter out the invalid trajectories. Particularly, the VLM examines the coherence and correctness of each constructed trajectory, and trajectories failing the check are discarded. Algorithm 1 details the labeling and filtering process for trajectory construction.

3.4 Video Application

With the processed videos, at inference time, the agent dynamically selects zero or one trajectory at each step to better decide the next action. We develop the following two-stage selection framework, which allows the agent to more accurately select a trajectory to focus on.

① *Stage 1: Coarse selection.* Given the current observation and task, the agent selects a candidate pool of trajectories from all videos of the task. Specifically, the VLM reads each trajectory’s objective and selects the top-3 relevant trajectories for each video. Since the input only contains the textual objectives without the observation and action sequence, the selection is efficient. ② *Stage 2: Detailed selection.* The agent then inspects the candidate pool more closely by additionally checking the initial observation and the action sequence of each trajectory. It selects a single trajectory that is most helpful for determining the next action,

Algorithm 1 Trajectory labeling and filtering

Input: action-observation sequence $\{(o_t, a_t)\}_{t=1}^N$

Output: set of labeled trajectories \mathcal{T}

$\mathcal{T} \leftarrow \emptyset$

for $i = 1$ **to** $N - 1$ **do**

for $j = i + 1$ **to** N **do**

 Use VLM to generate an objective $g_{i:j}$ for the action sequence $a_{i:j}$

if the sequence is complete and coherent **then**

$\mathcal{T} \leftarrow \mathcal{T} \cup \left\{ \left(g_{i:j}, \{(o_t, a_t)\}_{t=i}^j \right) \right\}$

end if

end for

end for

return \mathcal{T}

or no trajectory if all candidates are irrelevant. Finally, the selected trajectory, including the objective, observation, and action sequence, is provided in context when the agent generates the next action.

Maintaining coherent context across steps. Selecting trajectories at each step could lead to frequent context changes for the agent across consecutive steps, which might affect performance (*e.g.*, when the agent selects a trajectory and starts to follow, it could select a different trajectory at the next step, which deviates from the previous context and plan). To minimize unnecessary reselection, at each step, the agent first decides if the previously selected trajectory is still relevant and applicable for the current observation. If yes, we simply continue with the same trajectory. Otherwise, we repeat the two-stage selection to obtain a new trajectory. This keeps the in-context guidance stable when progress is steady, while allowing quick adaptation when the UI diverges.

4 Experiments

In this section, we conduct experiments to verify the effectiveness of our proposed method. Particularly, we aim to answer two questions: ① Does providing video tutorials improve agents’ performance on computer-use tasks? ② What factors affect the agents’ ability to learn from videos?

4.1 Experimental Setting

Evaluation benchmarks. We evaluate on two widely used benchmarks for computer-use tasks. ① **OSWORLD-VERIFIED** (Xie et al., 2024; 2025c) (hereafter as OSWORLD) contains 369 Ubuntu desktop tasks spanning real applications such as GIMP and VSCode, as well as tasks that involve multiple applications. Each task is associated with a reproducible initial state and manually designed evaluators. Following official guidance, we exclude 8 Google Drive tasks for ease of deployment. ② **WEBARENA** (Zhou et al., 2024) contains 812 tasks focusing on the web environment across different domains, such as GitLab and shopping. The agent needs to navigate the browser to complete a given task. As the official server for the map domain is not accessible, we evaluate on the remaining domains. For both benchmarks, we use the official evaluation implementation and report the final success rate (SR) over all tasks.

Baselines. We compare with three strong baselines. ① **Base agent:** We adopt the state-of-the-art agent frameworks on both benchmarks. These frameworks do not have access to external tutorials during inference. Specifically, on OSWORLD, we use **Jedi** (Xie et al., 2025b) as the base agent; on WEBARENA, we use **AgentOccam** (Yang et al., 2025b). ② **Text-only tutorials:** We retrieve online textual tutorials relevant to the task and prepend the retrieved tutorials to the context at each step. For retrieval, we follow Agent S (Agashe et al., 2025a) to first generate queries for each task and then use Perplexica Search Engine³ to retrieve and summarize online tutorials into a comprehensive document. ③ **Video transcript only:** We retrieve and filter online videos as described in Section 3.2, but we only use their transcripts without the visual information. To obtain the transcript, we use the **turbo** model from Whisper (Radford et al., 2023). Note that baselines ② and ③ are built upon the base agent frameworks in ① by providing external knowledge.

Agent settings. For all baselines and our method, we use **gpt-5-mini-2025-08-07** as the underlying VLM backbone at inference. On OSWORLD, all methods (ours and baselines) receive *only* the current screenshot as the observation, following prior work Xie et al. (2025b). We also set the maximum number of steps to 50. On WEBARENA, we follow Yang et al. (2025b) to use the accessibility tree of the page, as well as the additional current screenshot as the observation. The maximum steps is set to 20. Additional parameters of the agents are detailed in Appendix A.

Video sources. On OSWORLD, we directly retrieve videos from YouTube, as described in Section 3.2. After filtering, we collect videos for 211 out of 361 tasks,⁴ and on average, each task contains 3.6 videos. On WEBARENA, we use the videos from VideoWebArena (Jang et al., 2025a), which contains pre-recorded

³<https://github.com/ItzCrazyKns/Perplexica>.

⁴The remaining 150 tasks do not have a relevant video. On those tasks, our method behaves the same as the baseline agent.

Table 2: Success rate on WEBARENA.

	Shopping	Shopping Admin	Reddit	GitLab	All
AgentOccam (Yang et al., 2025b)	44.8	56.1	60.2	33.9	47.3
+ Transcript	44.6	54.2	53.6	42.0	47.9
Ours	47.1 (+2.3)	58.5 (+2.4)	60.6 (+0.4)	44.1 (+10.2)	51.5 (+4.2)

videos for 672 tasks in WEBARENA, excluding the tasks that require the map server. Note that although these videos are specifically designed to cover the skills required in WEBARENA, they do not align with the evaluation tasks exactly, and one video could contain multiple skills corresponding to many tasks. Thus, the agent still needs to select relevant trajectories from each video and adapt its content to tasks at hand, instead of directly copying the video.

Implementation details. We build our method based on Jedi and AgentOccam on the two benchmarks, respectively. For video retrieval and processing, we use Qwen2.5-VL-32B-Instruct (Bai et al., 2025) and Qwen3-30B-A3B-Instruct-2507 (Yang et al., 2025a) as the VLM and LLM, respectively. Additional implementation details are in Appendix A.

4.2 Main Results

Overall success rate. Table 1 shows the results on OSWORLD. We separately report the performance on the subset of 211 tasks for which we retrieve at least one relevant video. On the remaining 150 tasks, our method and the transcript baseline behave exactly the same as the Jedi framework. As the table indicates, our method outperforms all three baselines, especially on the subset that has video tutorials. This result confirms that our proposed pipeline improves the agent’s ability by learning from videos in context. Particularly, providing transcripts only does not improve the performance, which suggests that our structured, visually grounded demonstration trajectories are more useful than unstructured textual guidance during inference. Additionally, on average, our method takes 28.8 steps to complete a task, which is fewer than that of the Jedi baseline (30.3). This suggests that having access to external knowledge allows the agent to finish the tasks more directly without unnecessary exploration.

Table 1: Success rate on OSWORLD. “Has Videos” denotes the subset of 211 tasks for which we find at least one relevant video.

	Has Videos (211)	All (361)
Jedi (Xie et al., 2025b)	46.8	41.0
+ Text tutorial	44.0	39.0
+ Transcript	46.7	41.0
+ Video demo	50.3 (+3.5)	43.1 (+2.1)

Table 2 shows the results on WEBARENA across different domains. Since the websites are synthetically built and do not exist in the real world, we do not include the text tutorial baseline. As can be observed, the results show similar trends as in Table 1, where our method consistently achieves better performance than the baseline framework and the variant that adds transcript information. Particularly, in more challenging domains like GitLab, our method achieves a more significant improvement (+10.2), indicating a good utilization of the videos. By contrast, in easier domains like Reddit, the agent already has enough knowledge of how to complete the task (*e.g.*, how to create a post or make a comment), therefore adding additional information in the video does not make a big impact.

4.3 Additional Analyses

We now investigate how different factors in our pipeline affect the final performance. Specifically, we analyze the impacts from four aspects, including the number of videos per task, the proposed two-stage trajectory selection method, the filtering step in video processing, and the importance of visual information in the video. We conduct all the following analyses on the subset of 211 tasks in OSWORLD that have at least one retrieved video.

Number of videos per task. We first study how the number of videos per task affects the performance. Specifically, during inference, we restrict the agent’s access to only a subset of retrieved and filtered videos. Table 3 shows the performance of our method as we vary the amount of available videos. As can be seen, having access to more videos improves the overall performance, since the agent can select the demonstration trajectory from a larger pool. This suggests that **our method’s performance could potentially be scaled up when more relevant videos**

Table 3: Performance on OSWORLD with different numbers of videos per task.

	Success Rate
1 video per task	46.9
3.6 videos per task	50.3

The two-stage trajectory selection. We then explore the effect of our dynamic two-stage trajectory selection method, which allows the agent to focus on a local segment of the video at each step. Specifically, we compare with a variant of our method that disables this mechanism: instead of constructing trajectories for all subsequences of actions and dynamically selecting a trajectory at each step during inference, we always prepend the longest trajectory that corresponds to the entire video. If a task has

Table 4: Ablation study on OSWORLD.

	Success Rate
Ours	50.3
No trajectory selection	45.8
No action filtering	46.6
No visual information	45.6

multiple videos, we first prompt the agent to select the most relevant one and always provide that video during inference. The second row of Table 4 (no trajectory selection) shows that the performance of this variant drops significantly compared to our original method. This indicates the **importance of splitting the video into short trajectories and dynamically selecting suitable trajectories at each step.**

Action filtering in video processing. Next, we investigate the impact of filtering actions during video processing in Section 3.3. Recall that after labeling actions in the video, we filter actions and only keep those that are relevant to the main task of the video. To validate the impact of this filtering step, we compare with a variant of our method without this step, while keeping everything else the same, including the trajectory selection method. Results in Table 4 show that this variant (no action filtering) is worse than our original method, which illustrates the **benefits of filtering labeled actions.**

Visual information in trajectories. Finally, we study whether the visual information in trajectories contributes to the final performance. Specifically, we modify our method so that when providing the selected trajectory in context, we only provide the textual information (objective and action sequence), without the screenshot at each step. Table 4 shows that this variant (no visual information) is worse than our original method, which provides the sequence of screenshots in context. This indicates that **text-only summaries cannot cover all useful information in the video**, and adding visual information can further boost the performance of the agent.

5 Conclusion

We studied how computer-use agents can learn from online videos at inference time, inspired by how humans watch and imitate short segments of tutorials to acquire domain-specific procedural knowledge. We proposed a framework that retrieves and filters tutorial videos, converts them into structured, visually grounded demonstration trajectories using a VLM, and dynamically selects a single trajectory as in-context guidance at each step. Experiments show consistent gains over strong base agents and variants using only textual tutorials or transcripts, and analyses highlight the importance of trajectory segmentation and selection, action filtering, and visual information.

References

- Saaket Agashe, Jiuzhou Han, Shuyu Gan, Jiachen Yang, Ang Li, and Xin Eric Wang. Agent s: An open agentic framework that uses computers like a human. In *The Thirteenth International Conference on Learning Representations*, 2025a.
- Saaket Agashe, Kyle Wong, Vincent Tu, Jiachen Yang, Ang Li, and Xin Eric Wang. Agent s2: A compositional generalist-specialist framework for computer use agents, 2025b.
- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, 2017.
- Yusuf Aytar, Tobias Pfaff, David Budden, Thomas Paine, Ziyu Wang, and Nando de Freitas. Playing hard exploration games by watching youtube. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025.
- 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. Windows agent arena: Evaluating multi-modal os agents at scale, 2024.
- Dongping Chen, Yue Huang, Siyuan Wu, Jingyu Tang, Huichi Zhou, Qihui Zhang, Zhigang He, Yilin Bai, Chujie Gao, Liuyi Chen, Yiqiang Li, Chenlong Wang, Yue Yu, Tianshuo Zhou, Zhen Li, Yi Gui, Yao Wan, Pan Zhou, Jianfeng Gao, and Lichao Sun. GUI-world: A video benchmark and dataset for multimodal GUI-oriented understanding. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.
- Difei Gao, Lei Ji, Zechen Bai, Mingyu Ouyang, Peiran Li, Dongxing Mao, Qinchen Wu, Weichen Zhang, Peiyi Wang, Xiangwu Guo, Hengxu Wang, Luwei Zhou, and Mike Zheng Shou. Assistgui: Task-oriented desktop graphical user interface automation, 2024.
- Xiaoxiao Guo, Shiyu Chang, Mo Yu, Gerald Tesauro, and Murray Campbell. Hybrid reinforcement learning with expert state sequences, 2019.
- Lawrence Keunho Jang, Yinheng Li, Dan Zhao, Charles Ding, Justin Lin, Paul Pu Liang, Rogerio Bonatti, and Kazuhito Koishida. Videowebarena: Evaluating long context multimodal agents with video understanding web tasks. In *The Thirteenth International Conference on Learning Representations*, 2025a.
- Yunseok Jang, Yeda Song, Sungryull Sohn, Lajanugen Logeswaran, Tiange Luo, Dong-Ki Kim, Kyunghoon Bae, and Honglak Lee. Scalable video-to-dataset generation for cross-platform mobile agents, 2025b.
- Chengyou Jia, Minnan Luo, Zhuohang Dang, Qiushi Sun, Fangzhi Xu, Junlin Hu, Tianbao Xie, and Zhiyong Wu. AgentStore: Scalable integration of heterogeneous agents as specialized generalist computer assistant. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Findings of the Association for Computational Linguistics: ACL 2025*, Vienna, Austria, July 2025.
- Kevin Qinghong Lin, Linjie Li, Difei Gao, Qinchen WU, Mingyi Yan, Zhengyuan Yang, Lijuan Wang, and Mike Zheng Shou. VideoGUI: A benchmark for GUI automation from instructional videos. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.

- YuXuan Liu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Imitation from observation: Learning to imitate behaviors from raw video via context translation, 2018.
- Dunjie Lu, Yiheng Xu, Junli Wang, Haoyuan Wu, Xinyuan Wang, Zekun Wang, Junlin Yang, Hongjin Su, Jixuan Chen, Junda Chen, Yuchen Mao, Jingren Zhou, Junyang Lin, Binyuan Hui, and Tao Yu. Videoagenttrek: Computer use pretraining from unlabeled videos, 2025.
- OpenAI. Computer-using agent, 1 2025. URL <https://openai.com/index/computer-using-agent/>. Accessed: 2025-10-20.
- Tianyue Ou, Frank F. Xu, Aman Madaan, Jiarui Liu, Robert Lo, Abishek Sridhar, Sudipta Sengupta, Dan Roth, Graham Neubig, and Shuyan Zhou. Synatra: Turning indirect knowledge into direct demonstrations for digital agents at scale. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Anthropic PBC. Introducing computer use, a new claude 3.5 sonnet, and claude 3.5 haiku. <https://www.anthropic.com/news/3-5-models-and-computer-use>, Oct 2024. Accessed: 2025-10-20.
- Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao Li, Yunxin Li, Shijue Huang, Wanjuan Zhong, Kuanye Li, Jiale Yang, Yu Miao, Woyu Lin, Longxiang Liu, Xu Jiang, Qianli Ma, Jingyu Li, Xiaojun Xiao, Kai Cai, Chuang Li, Yaowei Zheng, Chaolin Jin, Chen Li, Xiao Zhou, Minchao Wang, Haoli Chen, Zhaojian Li, Haihua Yang, Haifeng Liu, Feng Lin, Tao Peng, Xin Liu, and Guang Shi. Ui-tars: Pioneering automated gui interaction with native agents, 2025.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine Mcleavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *Proceedings of the 40th International Conference on Machine Learning*. PMLR, 2023.
- Christopher Rawles, Sarah Clinckemaulle, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Marybeth Fair, Alice Li, William E Bishop, Wei Li, Folawiyi Campbell-Ajala, Daniel Kenji Toyama, Robert James Berry, Divya Tyamagundlu, Timothy P Lillicrap, and Oriana Riva. Androidworld: A dynamic benchmarking environment for autonomous agents. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Karl Schmeckpeper, Oleh Rybkin, Kostas Daniilidis, Sergey Levine, and Chelsea Finn. Reinforcement learning with videos: Combining offline observations with interaction. In Jens Kober, Fabio Ramos, and Claire Tomlin (eds.), *Proceedings of the 2020 Conference on Robot Learning*, volume 155 of *Proceedings of Machine Learning Research*, pp. 339–354. PMLR, 16–18 Nov 2021.
- Chan Hee Song, Yiwen Song, Palash Goyal, Yu Su, Oriana Riva, Hamid Palangi, and Tomas Pfister. Watch and learn: Learning to use computers from online videos, 2025.
- Hongjin Su, Ruoxi Sun, Jinsung Yoon, Pengcheng Yin, Tao Yu, and Sercan O Arik. Learn-by-interact: A data-centric framework for self-adaptive agents in realistic environments. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Zeyi Sun, Ziyu Liu, Yuhang Zang, Yuhang Cao, Xiaoyi Dong, Tong Wu, Dahua Lin, and Jiaqi Wang. Seagent: Self-evolving computer use agent with autonomous learning from experience, 2025.
- Weihao Tan, Wentao Zhang, Xinrun Xu, Haochong Xia, Ziluo Ding, Boyu Li, Bohan Zhou, Junpeng Yue, Jiechuan Jiang, Yewen Li, Ruyi An, Molei Qin, Chuqiao Zong, Longtao Zheng, YuJie Wu, Xiaoqiang Chai, Yifei Bi, Tianbao Xie, Pengjie Gu, Xiyun Li, Ceyao Zhang, Long Tian, Chaojie Wang, Xinrun Wang, Börje F. Karlsson, Bo An, Shuicheng YAN, and Zongqing Lu. Cradle: Empowering foundation agents towards general computer control. In *Forty-second International Conference on Machine Learning*, 2025.
- Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation, 2018.
- Faraz Torabi, Garrett Warnell, and Peter Stone. Generative adversarial imitation from observation, 2019.

- Haoming Wang, Haoyang Zou, Huatong Song, Jiazhan Feng, Junjie Fang, Junting Lu, Longxiang Liu, Qinyu Luo, Shihao Liang, Shijue Huang, Wanjun Zhong, Yining Ye, Yujia Qin, Yuwen Xiong, Yuxin Song, Zhiyong Wu, Aoyan Li, Bo Li, Chen Dun, Chong Liu, Daoguang Zan, Fuxing Leng, Hanbin Wang, Hao Yu, Haobin Chen, Hongyi Guo, Jing Su, Jingjia Huang, Kai Shen, Kaiyu Shi, Lin Yan, Peiyao Zhao, Pengfei Liu, Qinghao Ye, Renjie Zheng, Shulin Xin, Wayne Xin Zhao, Wen Heng, Wenhao Huang, Wenqian Wang, Xiaobo Qin, Yi Lin, Youbin Wu, Zehui Chen, Zihao Wang, Baoquan Zhong, Xinchun Zhang, Xujiang Li, Yuanfan Li, Zhongkai Zhao, Chengquan Jiang, Faming Wu, Haotian Zhou, Jinlin Pang, Li Han, Qi Liu, Qianli Ma, Siyao Liu, Songhua Cai, Wenqi Fu, Xin Liu, Yaohui Wang, Zhi Zhang, Bo Zhou, Guoliang Li, Jiajun Shi, Jiale Yang, Jie Tang, Li Li, Qihua Han, Taoran Lu, Woyu Lin, Xiaokang Tong, Xinyao Li, Yichi Zhang, Yu Miao, Zhengxuan Jiang, Zili Li, Ziyuan Zhao, Chenxin Li, Dehua Ma, Feng Lin, Ge Zhang, Haihua Yang, Hangyu Guo, Hongda Zhu, Jiaheng Liu, Junda Du, Kai Cai, Kuanye Li, Lichen Yuan, Meilan Han, Minchao Wang, Shuyue Guo, Tianhao Cheng, Xiaobo Ma, Xiaojun Xiao, Xiaolong Huang, Xinjie Chen, Yidi Du, Yilin Chen, Yiwen Wang, Zhaojian Li, Zhenzhu Yang, Zhiyuan Zeng, Chaolin Jin, Chen Li, Hao Chen, Haoli Chen, Jian Chen, Qinghao Zhao, and Guang Shi. Ui-tars-2 technical report: Advancing gui agent with multi-turn reinforcement learning, 2025a.
- Xinyuan Wang, Bowen Wang, Dunjie Lu, Junlin Yang, Tianbao Xie, Junli Wang, Jiaqi Deng, Xiaole Guo, Yiheng Xu, Chen Henry Wu, Zhennan Shen, Zhuokai Li, Ryan Li, Xiaochuan Li, Junda Chen, Boyuan Zheng, Peihang Li, Fangyu Lei, Ruisheng Cao, Yeqiao Fu, Dongchan Shin, Martin Shin, Jiarui Hu, Yuyan Wang, Jixuan Chen, Yuxiao Ye, Danyang Zhang, Dikang Du, Hao Hu, Huarong Chen, Zaida Zhou, Haotian Yao, Ziwei Chen, Qizheng Gu, Yipu Wang, Heng Wang, Diyi Yang, Victor Zhong, Flood Sung, Y. Charles, Zhilin Yang, and Tao Yu. Opencua: Open foundations for computer-use agents, 2025b.
- Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao Yu, and Lingpeng Kong. Os-copilot: Towards generalist computer agents with self-improvement, 2024.
- Jingxu Xie, Dylan Xu, Xuandong Zhao, and Dawn Song. Agentsynth: Scalable task generation for generalist computer-use agents, 2025a.
- 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. OSWorld: Benchmarking multimodal agents for open-ended tasks in real computer environments. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.
- Tianbao Xie, Jiaqi Deng, Xiaochuan Li, Junlin Yang, Haoyuan Wu, Jixuan Chen, Wenjing Hu, Xinyuan Wang, Yuhui Xu, Zekun Wang, Yiheng Xu, Junli Wang, Doyen Sahoo, Tao Yu, and Caiming Xiong. Scaling computer-use grounding via user interface decomposition and synthesis, 2025b.
- Tianbao Xie, Mengqi Yuan, Danyang Zhang, Xinzhuang Xiong, Zhennan Shen, Zilong Zhou, Xinyuan Wang, Yanxu Chen, Jiaqi Deng, Junda Chen, Bowen Wang, Haoyuan Wu, Jixuan Chen, Junli Wang, Dunjie Lu, Hao Hu, and Tao Yu. Introducing osworld-verified. *xlang.ai*, July 2025c. URL <https://xlang.ai/blog/osworld-verified>.
- Yiheng Xu, Dunjie Lu, Zhennan Shen, Junli Wang, Zekun Wang, Yuchen Mao, Caiming Xiong, and Tao Yu. Agenttrek: Agent trajectory synthesis via guiding replay with web tutorials. In *The Thirteenth International Conference on Learning Representations*, 2025a.
- Yiheng Xu, Zekun Wang, Junli Wang, Dunjie Lu, Tianbao Xie, Amrita Saha, Doyen Sahoo, Tao Yu, and Caiming Xiong. Aguis: Unified pure vision agents for autonomous GUI interaction. In *Forty-second International Conference on Machine Learning*, 2025b.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao

- Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025a.
- Ke Yang, Yao Liu, Sapana Chaudhary, Rasool Fakoor, Pratik Chaudhari, George Karypis, and Huzefa Rangwala. Agentoccam: A simple yet strong baseline for LLM-based web agents. In *The Thirteenth International Conference on Learning Representations*, 2025b.
- Yan Yang, Dongxu Li, Yutong Dai, Yuhao Yang, Ziyang Luo, Zirui Zhao, Zhiyuan Hu, Junzhe Huang, Amrita Saha, Zeyuan Chen, Ran Xu, Liyuan Pan, Silvio Savarese, Caiming Xiong, and Junnan Li. Gta1: Gui test-time scaling agent, 2025c.
- Tianhe Yu, Chelsea Finn, Annie Xie, Sudeep Dasari, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot imitation from observing humans via domain-adaptive meta-learning, 2018.
- Bofei Zhang, Zirui Shang, Zhi Gao, Wang Zhang, Rui Xie, Xiaojian Ma, Tao Yuan, Xinxiao Wu, Song-Chun Zhu, and Qing Li. Tongui: Building generalized gui agents by learning from multimodal web tutorials, 2025.
- Henry Hengyuan Zhao, Kaiming Yang, Wendi Yu, Difei Gao, and Mike Zheng Shou. Worldgui: An interactive benchmark for desktop gui automation from any starting point, 2025.
- Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. Gpt-4v(ision) is a generalist web agent, if grounded, 2024.
- Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web environment for building autonomous agents. In *The Twelfth International Conference on Learning Representations*, 2024.

Table 5: List of prompts in our pipeline.

		Prompt
Video Retrieval	Query generation	Figure 2
	Video filtering step 1	Figure 3
	Video filtering step 2	Figure 4
Video Processing	Action labeling	Figure 5
	Action merging	Figure 6
	Action filtering	Figure 7
	Trajectory objective generation	Figure 8
	Trajectory filtering	Figure 9
Video Application	OSWORLD Trajectory selection stage 1	Figure 10
	OSWORLD Trajectory selection stage 2	Figure 11
	WEBARENA Trajectory selection stage 1	Figure 12
	WEBARENA Trajectory selection stage 2	Figure 13
	Trajectory continuation	Figure 14

Table 6: Agent setting on OSWORLD.

	Value
Observation space	Screenshot
Maximum steps	50
Reasoning effort for VLM	medium
Temperature	Default

Table 7: Agent setting on WEBARENA.

	Value
Observation space	Screenshot & accessibility tree
Maximum steps	20
Reasoning effort for VLM	low
Temperature	Default

A Implementation Details

A.1 Video Retrieval and Processing

We follow the procedure described in Sections 3.2 and 3.3 to retrieve and process online videos. We use `Qwen2.5-VL-32B-Instruct` (Bai et al., 2025) and `Qwen3-30B-A3B-Instruct-2507` (Yang et al., 2025a) as the VLM and LLM, respectively. Table 5 summarizes the detailed prompts we use in our pipeline. Note that the video retrieval step is only done for OSWORLD. On WEBARENA, we use the pre-recorded videos from Jang et al. (2025a). Additionally, since each video in Jang et al. (2025a) covers knowledge for multiple tasks, we do not filter labeled actions on WEBARENA. Instead, we directly construct demonstration trajectories for all subsequences of actions.

A.2 Agent Framework

Our method is built on the `Jedi` (Xie et al., 2025b) and `AgentOccam` (Yang et al., 2025b) frameworks, respectively. We use `gpt-5-mini-2025-08-07` as the underlying VLM backbone for all methods. Table 5 shows the prompts we use for the trajectory selection at each step. Tables 6 and 7 summarize the agent parameters for OSWORLD and WEBARENA, respectively.

```

You are a search-query generator.
Given: (1) a laptop task a user wants to accomplish, and (2) the related desktop
applications.
Goal: Produce concise web search query that a typical user would type to find a tutorial or
guide to complete the task with those apps.

## Guidelines:
- Return only the plain query text (no labels, quotes, or operators like site:).
- Make it concise, use keywords and phrases, e.g., search for a specific application or a
specific website. Include less than 10 words in each query.
- **Abstract the general functionality or goal** behind the task. **Avoid task-specific
details** like filenames, personal data, file paths, page numbers, timestamps, etc.
- If the task involves multiple steps, focus on those that require the most application
knowledge.
- If the task involves only one app, generate exactly one query.
- If the task involves N apps (N > 1), you can generate one or more (but at most N) queries
, each at a new line. For example:
  - If there is a single primary app, generate a single query focused on that app.
  - If there are multiple primary apps, generate separate queries for each app.

```

Figure 2: Prompt used to generate search queries for video tutorials.

```

You are a helpful assistant.
Given: (1) a laptop task a user wants to accomplish, (2) the related desktop applications,
and (3) a list of retrieved videos (each with an ID, title, and description).
Goal: Select the videos that are most likely to help the user complete the task.

## Output Requirements
- Output a list of IDs for the selected videos. Output in JSON format with a single key `
selected_video_ids`. For example:

```json
{
 "selected_video_ids": [1, 5, 8, 12, 13, 17]
}
```

- Select at most 10 videos, ordered by relevance.
- If none of the videos are relevant, return an empty list.

## Guidelines:
- Exclude videos that cover unrelated applications or off-topic content, even if they are
tangentially related to the task.

```

Figure 3: Prompt used for coarse filtering of retrieved videos.

```

You are a helpful assistant.
Given: (1) a task a user wants to accomplish on their Ubuntu laptop, and (2) a retrieved
video (with title, description, transcript, and 10 sample frames).
Goal: Determine if the video (1) demonstrates the use of desktop applications on an Ubuntu
laptop, AND (2) can help the user complete the task.

## Guidelines:
- The video must contain directly applicable information that enables the user to complete
the task.

```

```

- The video must show actual usage of desktop applications (e.g., file managers, browsers,
terminals, IDEs, office software, media players, etc.) running on an Ubuntu operating
system.
- The demonstration should clearly show the interaction with the Ubuntu desktop interface (
e.g., mouse clicks, typing, application windows).
- A video should not be considered a demonstration if it only shows:
  - Slides or presentation decks
  - Talking-head segments (person speaking to the camera)
  - Non-desktop applications (e.g., mobile apps)

## Output Requirements
Generate your output in two sections:
OBSERVATIONS:
- List your observations about the video frames, focusing on elements that indicate whether
desktop applications are being used on an Ubuntu laptop.
- Be specific about what you see in the frames, such as application windows, desktop
environment features, or any other relevant details.

JUDGEMENT:
- Output your judgement in JSON format with a single key `judge`. For example:
```json
{"judge": <true or false>}
```
where `true` means the video is a demonstration of using Ubuntu desktop applications and
can help the user complete the task, and `false` means otherwise.

```

Figure 4: Prompt used for content verification of retrieved videos.

```

You are given consecutive frames from a screen recording. Based on the visual change
between these frames, infer the most likely user action(s) that caused the change.

Below is the list of possible actions:
- click [button]: Click on element button. E.g., `click the [Submit] button`.
- type [text] [box]: Type text in a box. E.g., `type [Hello World] in the [search box]`.
- right click [button]: Right click a button or a link. E.g., `right click the product link`.
- drag [element] to [position]: Drag the element to specified position. E.g., `drag [text
box] to top of the page`.
- press [key]: Press one or more keys. E.g., `press [Ctrl+F]` or `press [Esc]`.

## Instructions
1. Carefully analyze the visual changes between these frames to determine the most likely
user action(s) from the list above.
2. For each action, also indicate the start and end frame numbers (inclusive) where the
action was observed.

## Output Format
Generate your response in two sections:
OBSERVATION AND REASONING:
- Describe the visual changes step-by-step and explain what user action(s) likely caused
these changes.
ACTIONS:
- Conclude your response with all identified actions in JSON format. For example:
```json
[
 {"action": "click the [Submit] button", "start_frame": 1, "end_frame": 2},
 {"action": "type [Hello World] in the [search box]", "start_frame": 4, "end_frame": 7},
 // more actions if applicable
]
```

## Guidelines

```

```

- **Action Duration:** Include all frames from the **start of the action** (e.g., cursor moves onto the button) to the **end of the resulting change** (e.g., page finishes loading) .
- **Atomic Actions:** Break down actions into the smallest clear units. For example, two consecutive button clicks should be annotated as two separate `click` actions with different frame ranges. However, if a user types a sentence in a box, annotate it as a single typing action.
- **Evidence Requirement:** Only annotate a click or typing action when there is **clear, observable evidence** (e.g., a button is visibly pressed or text is entered). Ignore cursor movements or hovers without results.
- **Terminology Consistency:** Use the given action names. For instance, use "click" instead of "open" or "select".
- **No Overlaps:** Frame ranges for different actions should not overlap.
- **Empty Results:** If no action can be determined, return an empty list `[]`.

```

Figure 5: Prompt used for action labeling on OSWORLD. Note that we use the same prompt on WEBARENA except for the list of possible actions.

```

You are given a list of UI actions. Your task is to detect duplicate or equivalent actions.

Two actions are considered duplicates if they have the same action type and target the same UI element with the same intent.
If duplicates are found, merge them into a single action that accurately represents the combined intent.

## Output Format
Generate your response in two sections:
REASON:
- Describe the reasoning process for identifying duplicates.
MERGED ACTIONS:
- Conclude your response with all identified duplicates in JSON format. For example:
```json
[
 {"merged_action": "click the [Submit] button", "original_action_ids": [0, 4]}
 {"merged_action": "type [Hello World] in the [search box]", "original_action_ids": [2, 3, 5]}
]
```
where `merged_action` is the merged action, and `original_action_ids` is the list of original action IDs that were merged.

## Guidelines
- **Natural Continuation:** If two or more actions are partial segments of the same continuous action (e.g., two typing steps that form a single sentence), merge them into one complete action.
- **Do Not Merge Distinct Atomic Actions:** If actions have different purposes or target different UI elements, they must remain separate.
- **Preserve Action Type:** Keep the merged action in the same category as the original actions.
- **Empty Results:** If no duplicates are found, return an empty list.

```

Figure 6: Prompt used for merging duplicated actions.

You are given (1) information of a YouTube video (title, description, and transcript) and (2) a list of UI actions identified from the video. Your task is to select actions that are relevant to the task demonstrated in the video.

Output Format

Generate your response in two sections:

ANALYSES:

- Describe the task demonstrated in the video based on the title, description, and transcript.
- Analyze each action line by line, describe whether it is important or relevant to the task, and explain why.

KEPT ACTIONS:

- Output your results with a JSON list containing the indices of relevant actions. For example:

```
```json
[1, 2, 4, 6, 7, 8, 9, 12, 13, 17, 20]
```
```

Guidelines

- **Relevance to Task**: Remove actions that are tangential or unrelated to the primary task. Example of irrelevant actions include (but are not limited to):
 - Mouse movements or hover actions that do not lead to a click or interaction.
 - Waiting actions unless they are crucial for understanding the sequence of events.
 - Scrolling actions that do not contribute to the task.
- **Redundancy**: If multiple actions represent the same intent or outcome, keep the most comprehensive one and remove the others.

Figure 7: Prompt used to filter out irrelevant actions.

You will be given a screenshot of a {platform} before any actions, a sequence of actions performed on the {platform}, and a screenshot of the {platform} after the actions. Your task is to determine whether the sequence of actions is:

1. **Coherent**: The actions are possible given the initial screenshot and lead logically to the final screenshot.
2. **Complete**: All steps necessary to finish the task are included. If any essential step is missing, it is **not** complete.
3. **Contains no unnecessary actions**: Every action should be relevant to achieving the task. Extra, unrelated steps make the sequence invalid.
4. **Matching a specific, reasonable task**: The actions clearly and fully accomplish one identifiable goal.

If the sequence contains extra actions, is missing required steps, or does not fully achieve the task, output "No task".
 Example: If the intended task is creating a new issue, but there is no action clicking [Create issue], this is **not** a valid task.

Output Format

Generate your response in two sections:

OBSERVATION AND REASONING:

- Describe what is visible in the before and after screenshots.
- Go through the sequence of actions step-by-step, explaining their effects.
- State whether they accomplish a clear, reasonable task (and describe it).
- Confirm that all required steps for the task were performed - otherwise state that it is incomplete.
- Confirm that all steps are necessary and that there are no irrelevant actions.

TASK:

- If a clear, complete task is achieved, output:

```
```json
{"task": "<concrete description of the task>"}
```
```

```

'''
- If the actions are incomplete, incoherent, contain extra actions, or do not achieve any
reasonable task, output:
```json
{"task": "No task"}
```

## Example
```json
{"task": "Filter for issues labeled as bug."}
```
'''

```

Figure 8: Prompt used to generate objectives for a sequence of actions.

```

You will be given a task description and a trajectory that consists of a screenshot of a {
platform} before any actions, a sequence of actions performed on the {platform}, and a
screenshot of the {platform} after the actions. Your job is to decide whether the
trajectory successfully accomplishes the task according to the criteria below.

## Evaluation Criteria
1. Alignment: The trajectory's final state clearly fulfills the task described.
2. Coherence: Each action logically follows from the previous observation, and no
action contradicts the task.
3. Naturalness: The actions resemble realistic human interaction with the {platform} in
this context.
4. Reasonableness: The task is completed efficiently, without unnecessary steps,
backtracking, or overly complex methods.

## Output Format
Generate your response in two sections:
OBSERVATION AND REASONING:
- Describe what is visible in the before and after screenshots.
- Determine if the final state of the {platform} matches the expected outcome of the task.
- Go through the sequence of actions step-by-step, explaining their effects.
- Confirm that all necessary steps were performed and that each action was required.
JUDGMENT:
```json
{"judge": <true or false>, "reason": "<reasoning for the judgment>"}
```

```

Figure 9: Prompt used to filter labeled trajectories.

```

You are an autonomous intelligent agent tasked with completing a task on an Ubuntu desktop.
Given the progress so far, the current observation, and a list of tasks with human
demonstrations, your task is to select up to 3 most relevant demonstrations that can help
you determine the next action.

Output the id of the demonstration tasks inside triple backticks, split by commas. For
example, ```2, 15, 23```.

If you think none of the demonstrations are relevant, output ```None```.

```

Figure 10: Prompt used for stage 1 selection of demonstration trajectories on OSWORLD.

You are an autonomous intelligent agent tasked with completing a task on an Ubuntu desktop. Given the progress so far, the current observation, and a list of tasks with human demonstrations (each with a task description, the initial observation, and a sequence of actions), your task is to select the most relevant demonstration whose detailed information can help you determine the next action. Assume that for each demonstration, you may later access intermediate observations (i.e., after each action), which can further guide your decision-making.

Output the id of the demonstration inside triple backticks. For example, ```2``` or ```3```.

If you think none of the demonstrations are relevant, output ```None```.

Figure 11: Prompt used for stage 2 selection of demonstration trajectories on OSWORLD.

You are an autonomous intelligent agent tasked with navigating a web browser.

Given the progress so far, the current observation, and a list of demonstration tasks, your task is to select up to 3 most relevant demonstrations that can help determine the next action.

Output your response in three steps.

1. First, analyze whether you already have enough knowledge to determine the next action.
2. If not, inspect the demonstrations to see if they contain the required information.
3. Output your final decision in the following two formats:

Case 1 - no demonstration needed: If you already have enough knowledge to determine the next action, or none of the demonstrations are relevant, output None inside triple backticks. For example, ```None```.

Case 2 - demonstration(s) needed: If you need the detailed information from some demonstrations, output the id of the demonstration tasks inside triple backticks, split by commas. For example, ```2, 15```.

Figure 12: Prompt used for stage 1 selection of demonstration trajectories on WEBARENA.

You are an autonomous intelligent agent tasked with navigating a web browser.

Given the progress so far, the current observation, and a list of demonstration tasks extracted from a video (each with a task description, the initial observation, and a sequence of actions), your task is to select the most relevant demonstration that can help determine the next action.

Output your response in three steps.

1. First, analyze whether you already have enough knowledge to determine the next action.
2. If not, inspect the demonstrations to see if they contain the required information.
3. Output your final decision in the following two formats:

Case 1 - no demonstration needed: If you already have enough knowledge to determine the next action, or none of the demonstrations are relevant, output None inside triple backticks. For example, ```None```.

Case 2 - demonstration needed: If you need the detailed information from a specific demonstration, output the id of the demonstration inside triple backticks. For example, ```2``` or ```3```.

Figure 13: Prompt used for stage 2 selection of demonstration trajectories on WEBARENA.

```
You are an autonomous intelligent agent tasked with completing a task on an Ubuntu desktop.
You have previously found a demo trajectory for a related task, and now your task is to
determine if the demo is still relevant and useful for choosing your next action.

You will be given:
- The current state of your task and environment (i.e., your progress and current
  observation)
- The demo trajectory you found previously

Respond with one of the following, inside triple backticks:
* ```Yes``` - if the demo trajectory is still relevant and can help you determine the
  next action.
* ```No``` - if the demo trajectory is no longer relevant or helpful.
```

Figure 14: Prompt used to determine whether to preserve the previously selected trajectory.