
AgentStudio: A Toolkit for Building General Virtual Agents

Longtao Zheng^{1*}, Zhiyuan Huang^{3*}, Zhenghai Xue¹,
 Xinrun Wang¹, Bo An^{1,2}, Shuicheng Yan²
¹NTU, Singapore ²Skywork AI, Singapore ³ETH Zurich
<https://skyworkai.github.io/agent-studio/>

Abstract

Creating autonomous virtual agents capable of using arbitrary software on any digital device remains a major challenge for artificial intelligence. Two key obstacles hinder progress: insufficient infrastructure for building virtual agents in real-world environments, and the need for in-the-wild evaluation of fundamental agent abilities. To address this, we introduce AgentStudio, an online, realistic, and multimodal toolkit that covers the entire lifecycle of agent development. This includes environment setups, data collection, agent evaluation, and visualization. The observation and action spaces are highly generic, supporting both function calling and human-computer interfaces. This versatility is further enhanced by AgentStudio’s graphical user interfaces, which allow efficient development of datasets and benchmarks in real-world settings. To illustrate, we introduce a visual grounding dataset and a real-world benchmark suite, both created with our graphical interfaces. Furthermore, we present several actionable insights derived from AgentStudio, e.g., general visual grounding, open-ended tool creation, learning from videos, etc. We have open-sourced the environments, datasets, benchmarks, and interfaces to promote research towards developing general virtual agents for the future.

1 Introduction

Building autonomous virtual agents that can utilize every software tool on computers represents a longstanding goal in AI research (OpenAI, 2016). Such agents are designed to receive computer states (e.g., screenshots) and take actions through function calling or human-computer interfaces (e.g., keyboard and mouse) in response to natural language instructions. Fueled by the advancements in large language models (LLMs) (Brown et al., 2020; Chowdhery et al., 2022; OpenAI, 2023b, *inter alia*), there has been impressive progress towards general virtual agents, particularly in web (Kim et al., 2023; Zheng et al., 2023), desktop (Zhang et al., 2024), and video games (Tan et al., 2024; SIMA Team, 2024).

However, current research in virtual agents is hindered by two main challenges:

1. A lack of open and systematic infrastructure for both building and benchmarking agents in real-world computer control. i) Existing online interactive environments, such as locally hosted websites (Yao et al., 2022; Zhou et al., 2023), are restricted to domain-specific tasks and action spaces, thus not allowing for agent evaluation across the broad spectrum of human activities. Also, these environments lack features that allow for tool creation support, and natural language feedback. This can make it difficult for open-ended agents to self-correct and self-improve within context. ii) Static datasets (Deng et al., 2023; Rawles et al., 2023) are prone to be hacked and ignore multiple valid trajectories. Moreover, the dataset collection pipeline is mostly absent from existing datasets, which is essential for specializing agent abilities such as visual grounding.

*Equal contribution. This work was performed when Longtao Zheng and Zhenghai Xue were interns at Skywork AI.

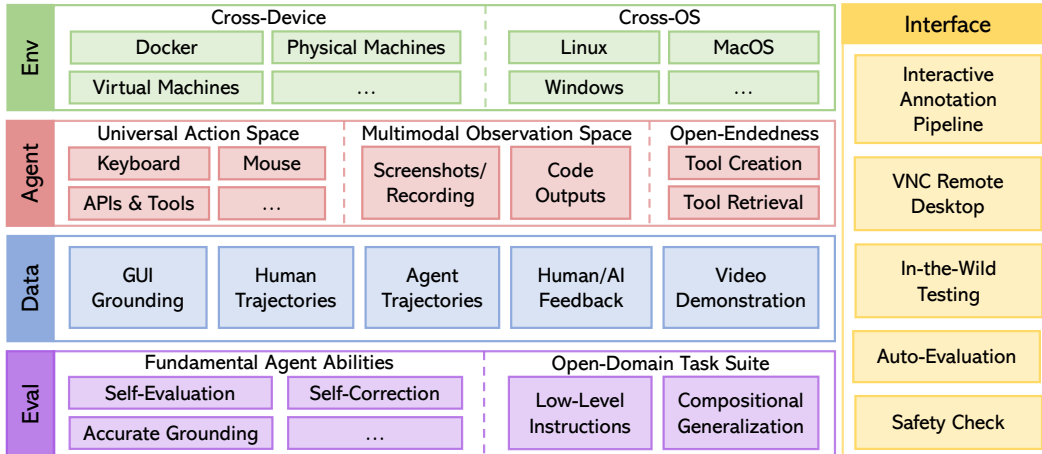


Figure 1: **The overview of AgentStudio framework.** We present a comprehensive toolkit covering the entire lifecycle of general virtual agents. The environment implementation is online, multimodal, and cross-platform. It offers unified observation and action spaces for agents, with native support for screen recording, open-ended tool use, etc. Furthermore, AgentStudio includes the complete pipeline for data annotation and in-the-wild evaluation, accompanied by interactive graphical user interfaces. This facilitates the creation of open-domain benchmarks and datasets that assess a wide range of agent capabilities.

2. The necessity for holistic evaluation of fundamental agent abilities in real-world scenarios. Current environments primarily focus on overall success rates instead of separately measuring core agent abilities, such as accurate graphical user interface (GUI) grounding, compositional generalization, self-evaluation as a critic, and utilizing both APIs or human-computer interfaces. Furthermore, they are not set up to test agents on open-domain and user-customized tasks based on human-computer interfaces, limiting how accurately their benchmark results can reflect real-world settings. However, these fundamental abilities and in-the-wild evaluation are critical for guiding the development of general virtual agents.

To address these issues, we release the public beta of AgentStudio¹, a toolkit encompassing the entire lifecycle of building and benchmarking general virtual agents *in the wild*. AgentStudio provides an integrated solution spanning environment setup, data collection, online testing, and result visualization (Figure 1). It adopts the following design choices:

The generic observation and action spaces consist of both human-computer interfaces and function calling. Specifically, agents interact with external environments through an interactive Python kernel. Therefore, they can automate keyboard-mouse interactions via Python scripts to control third-party applications, and leverage function calling when access to program internals or APIs is available. This universal action space enables agents to interact with arbitrary software. The observation space incorporates various modalities, ranging from screen recordings (videos) and screenshots (images) to code execution results (text), opening new research avenues such as learning from videos. Furthermore, AgentStudio **natively supports creating reusable code scripts as tools**. In contrast, previous environments used domain-specific observation and action spaces and did not consider open-ended tool creation, leading to limited generalizability of agents developed on them.

The environment implementation is online, realistic, and compatible with versatile operating systems and devices. This feature enables the development of agents that can handle massively open-domain and real-world scenarios, e.g., different system versions and resolutions. The online nature also facilitates research on agents that learn through trial-and-error interactions with environments. We also develop **two user-friendly graphical interfaces**: an interactive pipeline for collecting multimodal datasets of both human demonstrations and

¹We plan to expand the collection of environments, tasks, and data over time. Contributions and feedback from the community on how to improve this toolkit are welcome.

Environment	Domain	Online	Action Space	Image	Video	Human Feedback	Tool Creation	Data Pipeline
World of Bits (Shi et al., 2017)	Simplified Web	✓	Keyboard & Mouse	✓	✓	✗	✗	✗
AndroidEnv (Toyama et al., 2021)	Android	✓	Touchscreen	✓	✓	✗	✗	✗
WebGPT (Nakano et al., 2021)	Web-assisted QA	✓	Web Operations	✗	✗	✗	✗	✗
WebShop (Yao et al., 2022)	Web Shopping	✓	Web Operations	✓	✗	✗	✗	✗
Mind2Web (Deng et al., 2023)	Real-World Web	✗	Web Operations	✓	✗	✗	✗	✗
AITW (Rawles et al., 2023)	Android	✗	Touchscreen	✓	✗	✗	✗	✗
GALA (Mialon et al., 2023)	Tool-assisted QA	✗	APIs	✗	✗	✗	✗	✗
WebArena (Zhou et al., 2023)	Self-hosted Web	✓	Web Operations	✗	✗	✗	✗	✗
VisualWebArena (Koh et al., 2024)	Self-hosted Web	✓	Web Operations	✓	✗	✗	✗	✗
AgentStudio (Ours)	Real-World Devices	✓	Keyboard, Mouse and Code	✓	✓	✓	✓	✓

Table 1: **A comparison of environments for virtual agents.** AgentStudio is a platform that supports online interactions with versatile real-world computers, with the most generic observation and action spaces. It offers essential components for building general virtual agents that are missing from existing environments. More details can be found in Section 2.

agent synthetic trajectories, and a visualization interface for human-in-the-loop testing and gathering human feedback. Table 1 compares AgentStudio to other existing environments.

AgentStudio measures fundamental agent abilities and directs research towards improving these capabilities. We present two use cases to illustrate this design choice. Firstly, we utilize the AgentStudio interface to create a **GUI grounding dataset** across various applications and operating systems. This dataset consists of single-step atomic mouse operations with detailed instructions, e.g., clicking an icon. To complete these instructions, agents must generate accurate coordinates of mouse operations, a core ability that current benchmarks did not focus on. Using this dataset, we highlight the shortcomings in GUI grounding of current multimodal models. This also showcases the potential of scaling the grounding dataset with AgentStudio and training a model that can precisely translate single-step instructions into executable GUI actions. Secondly, we introduce an **real-world task suite** that ranges from simple function calling to complex cross-application tasks. Solving these tasks requires accurate GUI interaction, robust function calling, and long-horizon planning. AgentStudio accepts feedback from various sources, e.g., predefined rules, models, and humans. Additionally, we explore how current models perform in self-evaluation, i.e., serving as the critic for given trajectories. These two examples demonstrate how AgentStudio can facilitate building virtual agents and benchmarking in real-world scenarios.

Lastly, we highlight several promising research directions derived from AgentStudio, as well as the challenges we encountered during its development. The suggested research problems include auto-evaluation, visual grounding, etc. This showcases that AgentStudio has the potential to help future developments of general virtual agents.

In summary, this paper introduces an open and holistic toolkit for developing general virtual agents. We have open-sourced everything, including environment implementation, datasets, task suites, data collection pipeline, and graphical interfaces. We hope that AgentStudio will serve as a useful platform for the community to collectively scale data, develop novel algorithms, and benchmark agents with customized real-world tasks.

2 Related Work

AgentStudio builds on many previous efforts to build simulators and environments for automating computer operations using AI. For example, World of Bits (Shi et al., 2017) provided a minimalist web simulator, and AndroidEnv (Toyama et al., 2021) provided an emulator wrapper for Android devices. Early attempts in these environments primarily used deep reinforcement learning (Liu et al., 2018; Gur et al., 2018; Jia et al., 2018; Gur et al., 2021; Humphreys et al., 2022), which struggled to generalize to unseen tasks. Since these environments were tailored for reinforcement learning agents, they lack support for developing LLM-based agents with novel capabilities, e.g., code as policies and function calling (Liang et al., 2023). Additionally, with the progress made in LLMs, simplified

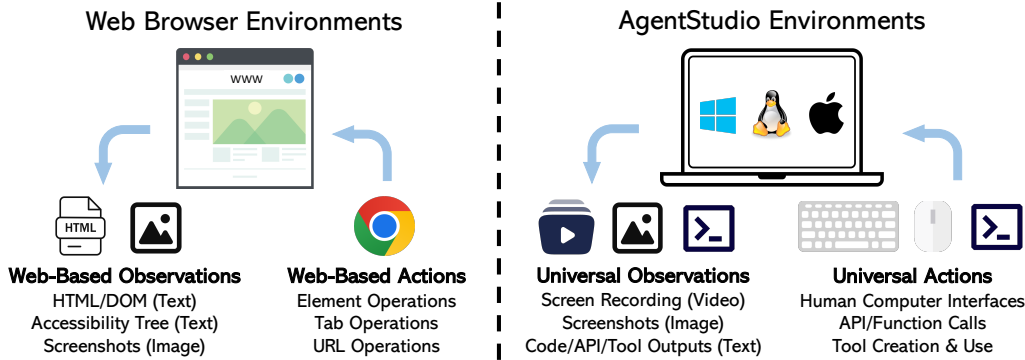


Figure 2: **A comparison between existing web browser environments and AgentStudio environments.** AgentStudio provides the most generic observation and action spaces, including human-computer interfaces and code execution. The implementation is also compatible with diverse operating systems and devices. These features drastically expand the potential task space and allow building and benchmarking agents in real-world settings.

environments such as World of Bits have been gradually solved (Gur et al., 2023b; Kim et al., 2023; Zheng et al., 2023; Gur et al., 2023a).

To build and benchmark virtual agents powered by LLMs in realistic and complex environments, WebGPT (Nakano et al., 2021) and WebShop (Yao et al., 2022) were introduced to measure web-assisted question answering and web shopping, respectively. However, they focused solely on web navigation, thus their tasks and action spaces are not generic enough for building general virtual agents. Similar issues apply to Mind2Web (Deng et al., 2023), WebArena (Zhou et al., 2023), and VisualWebArena (Koh et al., 2024). For example, the action space in Mind2Web only consists of clicking, selecting elements, and typing. Compared to AgentStudio, these environments were mainly designed to function as benchmarks rather than environments for agent development. As a result, they lack the essential infrastructure for building agents (Figure 1 and Table 1).

Unlike online environments, static datasets are vulnerable to hacking and cannot specify multiple different successful trajectories. More importantly, most works did not provide the pipeline for data collection. GAIA (Mialon et al., 2023) provided an easy-to-evaluate QA dataset requiring tool use to solve. AITW released a large-scale dataset of Android operations. ToolBench (Xu et al., 2023; Qin et al., 2023) and AgentBench (Liu et al., 2023) can only evaluate how well the agents have learned to use specific APIs, instead of general computer control. OpenAgents (Xie et al., 2023) focused on agent hosting and visualization while neglecting support for data collection, agent development, and evaluation. There are also concurrent works that introduce new static datasets to benchmark the agents they proposed, such as ScreenSpot (Cheng et al., 2024), ScreenAgent (Niu et al., 2024), OmniAct (Kapoor et al., 2024), WindowsBench (Zhang et al., 2024), and OS-Copilot (Wu et al., 2024). Similarly, an open data collection pipeline is missing in these works.

3 AgentStudio

AgentStudio is a holistic platform built on real-world environments to facilitate agent development. It extends beyond a typical benchmark with the following key features:

3.1 Universal Observation and Action Spaces

AgentStudio provides unified observation and action spaces that is compatible with both how humans interact with computers and function calling. This feature enables agent evaluation and data collection on any human-performed task, which drastically expands the potential task space. Therefore, AgentStudio can facilitate the development and benchmarking of agents that can generalize across diverse real-world use cases. In comparison,

most environments tailored the observation and action spaces to specific domains, such as web operations or API calls, as illustrated in Figure 2 and Table 1.

Action Space. The action space supports both high-level function calls and low-level atomic operations such as keyboard and mouse control. Similar to code interpreter (OpenAI, 2023a), agents can execute code via a Python kernel to call APIs, import tools, and control human-computer interfaces. Specifically, the control of human-computer interfaces (e.g., keyboard and mouse) is implemented through `pyautogui`, a Python package for GUI automation. This GUI control can be viewed as a form of tool use, and the implementation can be extended to support other interfaces such as touchscreens or other input modalities. This universal action space allows controlling arbitrary software. For example, agents can call APIs for applications with access to internals (e.g., shell, web search, and Google services). Even if there is no API access (e.g., closed-source third-party applications), agents can still automate operations via human-computer interfaces, the most general grounding.

Observation Space. The agent observations are multimodal, including videos of screen recordings, screenshots, and outputs from the Python kernel execution. The agent can further leverage tools or APIs to obtain structured internal states not visible to humans (e.g., HTML and accessibility tree). In comparison, most existing environments only provide a single screenshot or text observation after each action. This observation format avoids implementation difficulties such as specifying delay parameters after action execution. It also facilitates research on how agents can process complex and multimodal observations, and how they can solve tasks that require real-time video understanding.

Support for Tool Creation, Selection, and Use. AgentStudio allows open-ended agents to create and update tools by writing code into files, which can be reused later by importing them to the kernel. Therefore, agents can develop skills by combining basic operations or by creating tools to simplify the decision-making process, which is essential for building agents with compositional generalization and open-ended learning (Wang et al., 2023). In addition, AgentStudio can automatically convert code comments into tool documentation, which is designed to facilitate research on tool retrieval and documentation-based tool use.

3.2 Environments & Graphical Interfaces

AgentStudio aims to address the limitations of current benchmarks by targeting real-world scenarios with online interactions. We argue that the interactive nature of online environments is essential for research on exploration and self-correction of open-ended agents. To increase the accessibility of AgentStudio, we also develop an interactive annotation pipeline and a visualization interface for monitoring agent behaviors.

Online and Real-World Environments. While datasets are diverse and useful for training, they can be exploited, use domain-specific action space, and cannot recognize multiple valid solutions, rendering them ineffective for benchmarking. Meanwhile, most existing realistic online environments only have specific task suites and action spaces within limited domains. In contrast, AgentStudio adopts fully online interactive environments on real-world devices. It can operate in both local and remote mode. In remote mode, AgentStudio launches a VNC remote desktop, supporting all operating systems (e.g. Windows, Linux, and MacOS) and devices (e.g. Docker, virtual machines, and physical machines) that follow the VNC protocol. Both screen recording and action execution happen in *real time*, reflecting the complexity of real-world scenarios. Online and real-world environments allow agents to explore, learn from environment interactions, and accumulate new skills over time.

Natural Language Feedback and Visualization Interface. Compared to binary scalar rewards widely used in previous web environments and reinforcement learning environments, natural language feedback is considered more helpful for language agents to fix their own mistakes (Shinn et al., 2023). To promote research on in-context self-correction of virtual agents, AgentStudio supports generating natural language feedback from diverse sources, including rule-based evaluators, human feedback, and AI feedback. AgentStudio

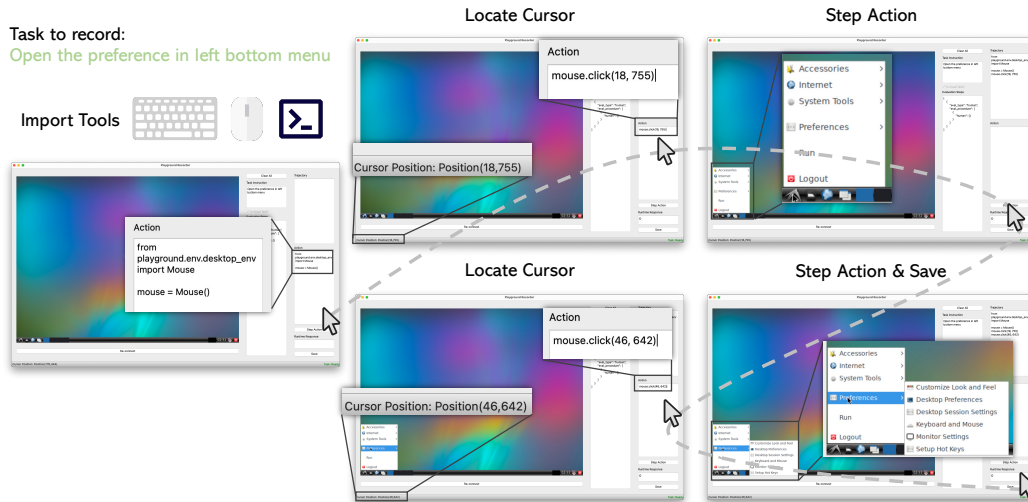


Figure 3: **Illustration of recording human annotation using AgentStudio.** This example showcases the process of recording complete keyboard-mouse action sequences required to open the preferences menu. The procedure involves the following steps: i) importing requisite tools, ii) obtaining the screen coordinates for mouse operations, and iii) executing actions. When collecting the GUI grounding dataset, users are additionally required to delineate the UI elements within a bounding box provided in the interface.

also offers a convenient visualization interface for monitoring agent behaviors, collecting human feedback, and testing customized tasks. In addition to in-context self-correction, the feedback gathered can be leveraged to finetune agents beyond human-level performance, similar to reinforcement learning from human feedback.

Interactive Data Collection Pipeline. Besides the visualization interface, AgentStudio provides an interactive pipeline for the collection of multimodal agent datasets (Figure 3). This is to bridge the gap that existing datasets did not release their data collection pipeline (Deng et al., 2023; Rawles et al., 2023). It allows us to create large-scale or task-specific datasets to train agents. For example, users can collect a dataset of a specific application to enhance agent performance. We believe that scaling data is the most effective approach to achieve accurate low-level GUI control, and AgentStudio offers a user-friendly pipeline for this.

4 Building Datasets and Benchmarks with AgentStudio

In this section, we present two downstream applications: i) utilizing the interactive annotation pipeline to collect GUI grounding data and evaluate current multimodal models, and ii) using the visualization interface for real-world benchmarking with customized task suites. These case studies demonstrate that AgentStudio offers a multimodal, simple-to-evaluate, and expandable pipeline for measuring and training agents, and sidesteps the issues of current agent evaluation, e.g., unrealistic, ignoring core abilities, and vulnerable to hacks.

4.1 GUI Grounding Dataset

GUI grounding with accurate coordinates is the main challenge for current virtual agents, since the interactable elements are not readily available as in web browsers (Koh et al., 2024). It has also been validated that state-of-the-art models can generate correct textual descriptions of the action, but struggle to ground them into environment operations (Zheng et al., 2024b). Therefore, we collect a dataset of tasks that only require a single step to complete, with clear instructions, e.g. clicking a specific button. This level of tasks mainly tests the visual grounding abilities of agents. Specifically, we define each data entry as a

3-tuple, $T = (g, s, a)$, where g is a detailed low-level instruction, s is the screenshot, and a is the mouse action. For example, given a screenshot s of a music player interface and the instruction g , “Click the Play button”, the action a is annotated as a bounding box of the button. Using AgentStudio, we devise a dataset with 227 instructions for mouse clicks, covering three popular desktop operating systems and nine applications.

Model	Windows				Linux			MacOS	
	OS	Games	PowerPoint	Word	OS	Browser	Calculator	Music	iMovie
Qwen-VL-Chat	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
SeeClick	42.3	47.6	24.0	9.1	60.7	9.1	69.0	53.8	35.7
Gemini-1.0 Pro	7.7	9.5	8.0	4.5	10.7	0.0	0.0	3.8	0.0
Claude-3 Sonnet	3.8	4.8	24.0	9.1	7.1	4.5	0.0	30.8	0.0
GPT-4V (1106)	11.5	23.8	20.0	9.1	14.3	4.5	3.4	11.5	3.6

Table 2: Experiment results on the GUI grounding dataset.

Results. Based on the dataset, we evaluate the GUI grounding abilities of current multi-modal models across various software and operating systems (Table 2). For closed-source APIs, we compare Gemini-1.0 Pro, Claude-3 Sonnet, and GPT-4V (1106). Gemini shows some ability to interact with or recognize GUI elements, especially on Windows, but fails for most tasks on other operating systems. This suggests the model may not generalize well across different operating systems or that the training data may be biased. Similarly, Claude-3 Sonnet scores highest on Windows PowerPoint and Apple Music but performs poorly on other applications. In contrast, GPT-4V demonstrates the highest overall grounding scores across most categories. However, it still falls short of the robustness required for real-world deployment, indicating the need for more data to improve the GUI grounding abilities of current multimodal models. For open-source models, we choose Qwen-VL and SeeClick (a Qwen-VL variant fine-tuned on existing GUI datasets (Zheng et al., 2023; Deng et al., 2023; Rawles et al., 2023)). Qwen-VL fails to obtain any successful outcomes, while SeeClick consistently achieves the highest scores in various settings. This highlights the importance of further scaling GUI grounding data to improve multimodal models in a data-driven approach for effective real-world deployment. Notably, SeeClick is a specialized model designed to output only coordinates instead of code. Therefore, in interpreting SeeClick’s performance in Table 2, it is presumed that all clicks are accurately matched.

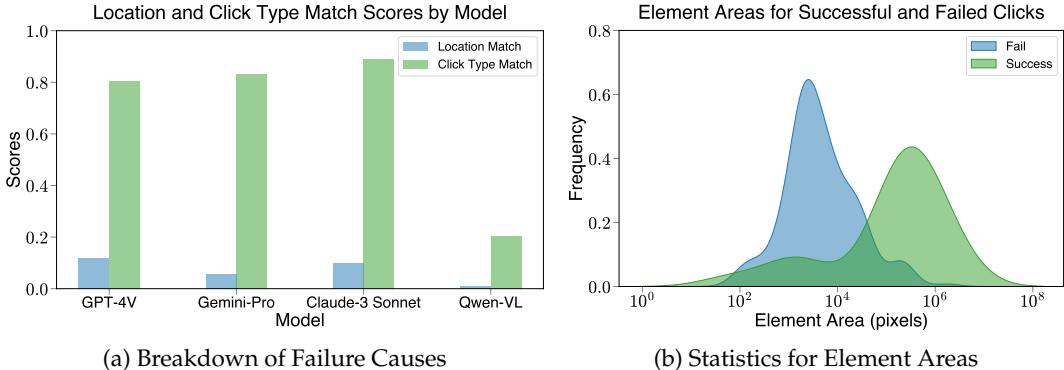


Figure 4: Comparative analysis of GUI grounding success rates.

Analysis. Figure 4a compares the failure causes based on two metrics: location match and click type match. These metrics measure the accuracy of the models predicting the GUI element location and the click type to execute (e.g., single click, double click, and right click). All models show low scores in precisely locating mouse action in the element, indicating much room for improvement. Figure 4b suggests that there is a correlation between the

size of GUI elements and the success rate of mouse operations. This implies that using a divide-and-conquer approach by splitting the screen into multiple areas and specifying coordinates within these smaller regions could probably improve GUI grounding.

4.2 Real-World Cross-Application Benchmark Suite

Unlike single-step GUI grounding, static datasets are less effective for benchmarking multi-step tasks because they cannot accept multiple valid solutions. Meanwhile, most existing online benchmarks are domain-specific and difficult to extend. To illustrate how AgentStudio facilitates agent evaluation in complex and real-world scenarios, we introduce a benchmark suite consisting of 77 real-world tasks. Despite being conceptually simple, these tasks pose challenges for current state-of-the-art virtual agents. Solving tasks in this suite requires various fundamental agent abilities, such as tool use, GUI grounding, compositional generalization, and long-horizon planning.

A task is formalized as a tuple: $T = (g, f_{\mathcal{R}}, f_{\mathcal{E}})$, where g is a natural language description of the task instruction, $f_{\mathcal{R}}$ optionally resets the environment, and $f_{\mathcal{E}}$ optionally evaluates the outcome trajectory. For example, for the instruction “Create a one-hour event Team Meeting in Google Calendar at 5pm today”, $f_{\mathcal{R}}$ checks and removes any existing events with the same name, while $f_{\mathcal{E}}$ evaluates success by checking via API if the event was properly created after the agent’s actions. The task suite is divided into three levels. Level 1 consists of 19 simple file manipulation tasks that can be accomplished by basic functions calling. Level 2 comprises 40 tasks involving Google services such as Gmail, Google Docs, Google Calendar, etc. This level represents common daily scenarios for virtual agents. These tasks can also be completed via API calls, though the APIs are more complicated than those in level 1. Level 3 covers 18 highly challenging tasks that require both GUI grounding and long-horizon decision-making. For example, agents may need to operate across multiple applications like Visual Studio Code and a video player to finish a task. These tasks are mostly compositional, cross-application, or involve complex GUI operations. Various automatic evaluation protocols are supported. Level 1 and 2 tasks are evaluated through a rule-based automatic evaluator, allowing for convenient benchmarking. For level 3 tasks, we rely on the human feedback interface in AgentStudio for evaluation.

Model	File Manipulation		Google Service		GUI & Cross-Application	
	Success	Critic Accuracy	Success	Critic Accuracy	Success	Critic Accuracy
GPT-3.5 Turbo	57.9	84.2	57.9	62.5	-	-
Gemini-1.0 Pro	78.9	89.5	40.0	72.5	11.1	88.9
GPT-4	100.0	52.6	65.0	47.5	0.0	100.0

Table 3: Experiment results on the real-world benchmark suite.

Results & Analysis. We compare the performance of three models in Table 3, including GPT-3.5 Turbo, Gemini-1.0 Pro, and GPT-4. For GUI & Cross-Application tasks, we use the vision versions of these models. Each task is evaluated based on two metrics: the percentage of tasks the model successfully completes (Success), and the accuracy in evaluating the success of trajectories (Critic Accuracy). The results suggest that while GPT-4 excels in most API calling tasks, it faces challenges in GUI and compositional tasks. Furthermore, GPT-4 has a lower critic accuracy compared to the other two models in level 1 and 2 tasks. On the other hand, Gemini-1.0 Pro and GPT-3.5 Turbo exhibit relatively lower success rates overall, but their higher critic accuracy implies that these models may have the potential to improve their performance by developing novel self-correcting algorithms.

5 Actionable Insights

AgentStudio has the potential to promote several lines of research that are essential for building general virtual agents. Here, we offer several possible research problems that are emphasized in the design choices of AgentStudio or encountered during its development.

General GUI Grounding. It has been validated that GPT-4V can achieve reasonable performance if there exists an oracle grounding (Zheng et al., 2024a), where the oracle can convert natural language instructions into correct GUI actions. Therefore, it is necessary to either train a specialized low-level visual grounding model, or develop novel prompting techniques that can accurately translate clear instructions into executable actions. Similar to PaLM-E (Driess et al., 2023), virtual agents may need LLMs to decompose tasks into detailed instructions on text space, and then ground instructions into executable GUI actions. To accomplish this goal, the data collection pipeline in AgentStudio can be used to collect grounding datasets for finetuning and benchmarking purposes. Additionally, the visualization interface in AgentStudio can be used to identify and analyze failure cases.

Learning from Documents and Video Demonstrations. Internet-scale data have contributed to remarkable success in language modeling. Similar to MineCLIP (Fan et al., 2022) and VPT (Baker et al., 2022), which leverages Internet-scale Minecraft data, the development of virtual agents can also benefit from the wealth of videos and documents available on the web. Furthermore, AgentStudio’s focus on video recording also enables the possibility of collecting video demonstrations labeled with actions. In addition to imitation learning, these video demonstrations can be used to train inverse dynamics models (Baker et al., 2022). Such models can add pseudo action labels to video recordings, which can potentially help leverage unsupervised internet-scale video data for training virtual agents.

Tool Creation, Selection, and Use. AgentStudio provides implementations that allow virtual agents to readily create and use tools. These capabilities are crucial for applying open-ended agent frameworks, e.g., Voyager (Wang et al., 2023), into real-world digital worlds. Tool creation and use can significantly reduce the compounded error in sequential decision-making and leverage knowledge learned from prior environment interactions. For example, when encountering unseen applications, an ideal virtual agent can explore the environment and write a program to use the application. This program will be added to a library of skills or tools, which can be reused. This serves as a “gradient-free” learning procedure, allowing the agent to adapt without updating its model parameters. AgentStudio turns the comments within the tool library into tool documentation, allowing tool use in a zero- or few-shot manner by reading documentation (Hsieh et al., 2023). It is also worthwhile to explore methods for selecting relevant tools from a large set of candidates.

A Generalist Critic Model. During benchmark creation, we found it challenging to provide feedback for open-domain tasks. Previous environments (Zhou et al., 2023) examined success rates through hand-written rules such as string matching. Similarly, AgentStudio extends functional evaluation to more domains and applications with its flexible evaluator framework. However, human-written evaluation is task-specific, not scalable, and cannot incorporate human preferences. Therefore, a general critic model that judges task completion and provides feedback is beneficial. It can auto-evaluate through API calling or judge based on trajectories. If a task fails, it provides natural language feedback to facilitate reflection and self-correction. With this critic, we can also leverage reinforcement learning from feedback to align with user preferences, which can be collected through AgentStudio.

6 Conclusion

In this work, we introduce AgentStudio, an open toolkit for developing general-purpose agents capable of operating in digital worlds. AgentStudio provides environment implementations in real-world settings and offers a holistic toolkit including data collection evaluation, visualization, and user interfaces. It enables developing and testing agents on arbitrary human-performed tasks. To illustrate the applications of AgentStudio, we have released a dataset focused on visual GUI grounding and a real-world benchmark suite with auto-evaluation and human evaluation. We acknowledge that our work has certain limitations. For example, the real-world complexity significantly increases the difficulty of checking functional correctness automatically. Therefore, we present several actionable insights derived from our framework, highlighting how AgentStudio can catalyze research efforts towards building versatile AI agents for digital worlds.

Reproducibility Statement

To ensure reproducibility, all resources such as code, datasets, agent trajectories, and leaderboard have been made publicly available at <https://skyworkai.github.io/agent-studio/>.

Ethic Statement

Autonomous agents that can control real-world digital devices have inherent risks, such as deleting important files and sending spam emails. To mitigate these risks, our framework requires user examination and confirmation before each code execution. Users should take extra caution when using these agents to interact with real-world environments.

References

- Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampedro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos. *Advances in Neural Information Processing Systems*, 35: 24639–24654, 2022.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901, 2020.
- Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. SeeClick: Harnessing GUI grounding for advanced visual GUI agents. *arXiv preprint arXiv:2401.10935*, 2024.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- 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*, 2023.
- Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems*, 35:18343–18362, 2022.
- Izzeddin Gur, Ulrich Rueckert, Aleksandra Faust, and Dilek Hakkani-Tur. Learning to navigate the web. In *International Conference on Learning Representations*, 2018.
- Izzeddin Gur, Natasha Jaques, Yingjie Miao, Jongwook Choi, Manoj Tiwari, Honglak Lee, and Aleksandra Faust. Environment generation for zero-shot compositional reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. A real-world WebAgent with planning, long context understanding, and program synthesis. In *The Twelfth International Conference on Learning Representations*, 2023a.
- Izzeddin Gur, Ofir Nachum, Yingjie Miao, Mustafa Safdari, Austin V Huang, Aakanksha Chowdhery, Sharan Narang, Noah Fiedel, and Aleksandra Faust. Understanding HTML with large language models. In *ICLR 2023 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2023b.

-
- Cheng-Yu Hsieh, Si-An Chen, Chun-Liang Li, Yasuhisa Fujii, Alexander Ratner, Chen-Yu Lee, Ranjay Krishna, and Tomas Pfister. Tool documentation enables zero-shot tool-usage with large language models. *arXiv preprint arXiv:2308.00675*, 2023.
- Peter C Humphreys, David Raposo, Tobias Pohlen, Gregory Thornton, Rachita Chhaparia, Alistair Muldal, Josh Abramson, Petko Georgiev, Adam Santoro, and Timothy Lillicrap. A data-driven approach for learning to control computers. In *International Conference on Machine Learning*, pp. 9466–9482. PMLR, 2022.
- Sheng Jia, Jamie Ryan Kiros, and Jimmy Ba. DOM-Q-NET: Grounded RL on structured language. In *International Conference on Learning Representations*, 2018.
- Raghav Kapoor, Yash Parag Butala, Melisa Russak, Jing Yu Koh, Kiran Kamble, Waseem Alshikh, and Ruslan Salakhutdinov. Omniact: A dataset and benchmark for enabling multi-modal generalist autonomous agents for desktop and web. *arXiv preprint arXiv:2402.17553*, 2024.
- Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language models can solve computer tasks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks. *arXiv preprint arXiv:2401.13649*, 2024.
- Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 9493–9500. IEEE, 2023.
- Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. Reinforcement learning on web interfaces using workflow-guided exploration. In *International Conference on Learning Representations*, 2018.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. AgentBench: Evaluating LLMs as agents. In *The Twelfth International Conference on Learning Representations*, 2023.
- Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. GAIA: A benchmark for general AI assistants. In *The Twelfth International Conference on Learning Representations*, 2023.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. WebGPT: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- Runliang Niu, Jindong Li, Shiqi Wang, Yali Fu, Xiyu Hu, Xueyuan Leng, He Kong, Yi Chang, and Qi Wang. Screenagent: A vision language model-driven computer control agent. *arXiv preprint arXiv:2402.07945*, 2024.
- OpenAI. Universe. <https://openai.com/blog/universe/>, 2016.
- OpenAI. Function calling and other api updates. <https://openai.com/blog/function-calling-and-other-api-updates>, 2023a.
- OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023b.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*, 2023.
- Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. Android in the wild: A large-scale dataset for Android device control. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.

-
- Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Percy Liang. World of Bits: An open-domain platform for web-based agents. In *International Conference on Machine Learning*, pp. 3135–3144. PMLR, 2017.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- SIMA Team. Scaling instructable agents across many simulated worlds. 2024.
- Weihao Tan, Ziluo Ding, Wentao Zhang, Boyu Li, Bohan Zhou, Junpeng Yue, Haochong Xia, Jiechuan Jiang, Longtao Zheng, Xinrun Xu, et al. Towards general computer control: A multimodal agent for red dead redemption ii as a case study. *arXiv preprint arXiv:2403.03186*, 2024.
- Daniel Toyama, Philippe Hamel, Anita Gergely, Gheorghe Comanici, Amelia Glaese, Zafarali Ahmed, Tyler Jackson, Shibl Mourad, and Doina Precup. AndroidEnv: A reinforcement learning platform for Android. *arXiv preprint arXiv:2105.13231*, 2021.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023.
- Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao Yu, and Lingpeng Kong. Os-Copilot: Towards generalist computer agents with self-improvement. *arXiv preprint arXiv:2402.07456*, 2024.
- Tianbao Xie, Fan Zhou, Zhoujun Cheng, Peng Shi, Luoxuan Weng, Yitao Liu, Toh Jing Hua, Junning Zhao, Qian Liu, Che Liu, et al. OpenAgents: An open platform for language agents in the wild. *arXiv preprint arXiv:2310.10634*, 2023.
- Qiantong Xu, Fenglu Hong, Bo Li, Changran Hu, Zhengyu Chen, and Jian Zhang. On the tool manipulation capability of open-source large language models. *arXiv preprint arXiv:2305.16504*, 2023.
- Shunyu Yao, Howard Chen, John Yang, and Karthik R Narasimhan. WebShop: Towards scalable real-world web interaction with grounded language agents. In *Advances in Neural Information Processing Systems*, 2022.
- Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang, Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qingwei Lin, Saravan Rajmohan, et al. Ufo: A ui-focused agent for windows os interaction. *arXiv preprint arXiv:2402.07939*, 2024.
- Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. Gpt-4v (ision) is a generalist web agent, if grounded. *arXiv preprint arXiv:2401.01614*, 2024a.
- Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. Gpt-4v(ision) is a generalist web agent, if grounded. *arXiv preprint arXiv:2401.01614*, 2024b.
- Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. Synapse: Trajectory-as-exemplar prompting with memory for computer control. In *The Twelfth International Conference on Learning Representations*, 2023.
- Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, et al. WebArena: A realistic web environment for building autonomous agents. In *The Twelfth International Conference on Learning Representations*, 2023.