

WEBCANVAS: BENCHMARKING WEB AGENTS IN ONLINE ENVIRONMENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

For web agents to be practically useful, they must adapt to the continuously evolving web environment characterized by frequent updates to user interfaces and content. However, most existing benchmarks only capture the *static* aspects of the web. To bridge this gap, we introduce *WebCanvas*, an innovative online evaluation framework for web agents that effectively addresses the dynamic nature of web interactions. *WebCanvas* contains three main components to facilitate realistic assessments: (1) A novel evaluation metric which reliably capture critical intermediate actions or states necessary for task completions while disregarding noise caused by insignificant events or changed web-elements. (2) A benchmark dataset called *Mind2Web-Live*, a refined version of original *Mind2Web* static dataset containing 542 tasks with 2439 intermediate evaluation states; (3) Lightweight and generalizable annotation tools and maintenance pipelines that enables the community to collect and maintain the high-quality, up-to-date dataset. Building on *WebCanvas*, we open-source a baseline agent framework with extensible modules for reasoning, providing a foundation for the community to conduct online inference and evaluations. Our best-performing agent achieves a task success rate of 23.1% and a task completion rate of 48.8% on the *Mind2Web-Live* test set. Additionally, we analyze the performance discrepancies across various websites, domains, and experimental environments. We encourage the community to contribute further insights on online agent evaluation, thereby advancing this field of research.

1 INTRODUCTION

The enhanced reasoning capabilities of foundational models (Ouyang et al., 2022; Achiam et al., 2023; Touvron et al., 2023a;b; Liu et al., 2024a; Bai et al., 2023) demonstrate the potential for autonomous agents performing on navigation and information retrieval tasks in real-time within web environment, thereby augmenting the human workforce (Shi et al., 2017; Nakano et al., 2021). However, the journey towards autonomous web agents delivering accurate, robust, fast, and cost-effective outcomes to end-users remains fraught with challenges. These include the inherent scarcity of data, the lack of knowledge and reasoning abilities of high-level actions on certain websites, and the absence of accurate and effective process feedback during execution, among others (Gür et al., 2023; Gur et al., 2024). We posit that a significant barrier to realizing the value of web agents is the lack of an easy-to-use platform for the community to drive effort towards real-time data gathering and web agent online benchmarking. This belief is grounded in following observations.

Digital agents require environmental observations and feedback for context. Thus, dynamic, real-world environments are essential for agent evaluation and data collection. The Internet itself emerges as the most extensive arena for the assessment of agents, offering an unparalleled complexity for environmental interaction (Liu et al., 2018; Zhou et al., 2023). However, the rapid evolution of the web environment introduces significant data distribution shifts over time. Figure 2 summarizes three prevalent patterns of changes in web tasks over time. For example, the *Mind2Web* dataset (Deng et al., 2024), which archives web-based interactions as static HTML snapshots and was released one year ago, shows that more than half tasks (50.5%) changed to some extent in their corresponding live websites one year later. This shift may potentially create discrepancies between the offline and online development and evaluation of real-world web agents. In addition, the accumulated knowledge and training data of static websites leads to the saturation of existing benchmarks, making

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

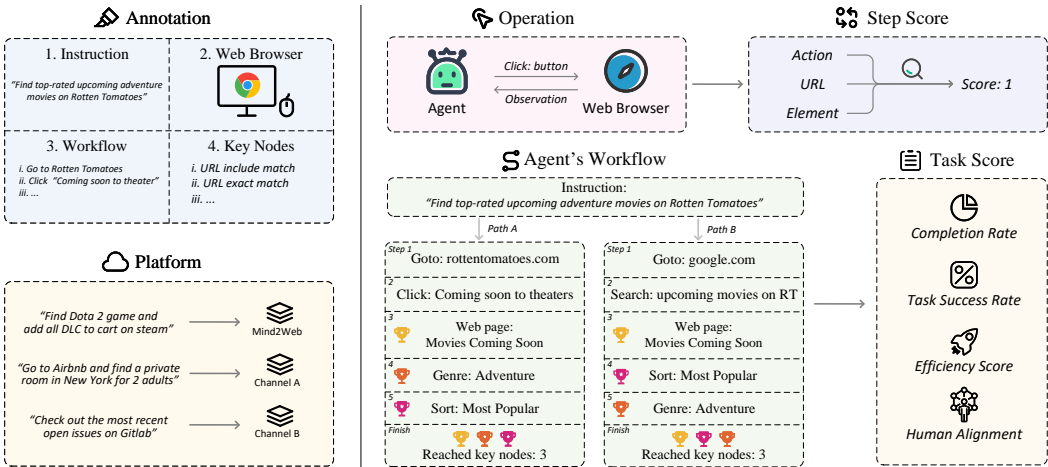


Figure 1: WebCanvas framework. The left side depicts the annotation process addressing each task, while the right side demonstrates the evaluation process during inference time, which involves collection of predicted actions, URLs, and elements targeted for interaction in online web environment, allowing for dynamic assessment. The framework accounts for the non-uniqueness of paths in online web interactions, with “Trophies” representing step scores earned upon successfully reaching each key node. Data maintenance pipeline of WebCanvas is detailed in §4.2.

it increasingly difficult to compare models and reasoning frameworks fairly and rigorously. We found the MindAct model trained in 2023 outperformed closely-held models like GPT-3.5 (Ouyang et al., 2022) and GPT-4 (Achiam et al., 2023) in Mind2Web static test set, but lagged behind in 2024 online evaluations (§5.1). Although previous works have attempted to evaluate the performance of web agents in online environments through human assessments (Zheng et al., 2024; He et al., 2024), achieving an objective, quantitative, and reproducible evaluation remains challenging.

To bridge this gap, we introduce WebCanvas, a dynamic and real-time framework designed for online evaluation of web agents with three key features. (1) **Progress-aware evaluation with key node annotation.** Existing evaluation metrics that focuses on action prediction accuracy (Deng et al., 2024; Zheng et al., 2024) can falsely penalize valid alternative solutions while outcome-based evaluation (Zhou et al., 2023; Koh et al., 2024; Mialon et al., 2023) requires fully reproducible standalone web environments. To address this gap, we introduce a novel concept termed “key nodes” – essential milestones that any task process must traverse, irrespective of the path taken. A comparative illustration of key nodes with these existing methods is provided in Figure 3. Key node annotation allows for a detailed, continuous analysis of agent behaviors, thereby enhancing insight into their decision-making strengths and weaknesses. (2) **Collaborative platform for community-driven annotations.** WebCanvas supports recording and annotation of web-based tasks and their corresponding key node evaluation through an advanced recording browser plugin with transparent data access. Furthermore, we have open-sourced an agent reasoning framework that enhances the integration and customization of various agent modules for online web tasks. This initiative provided guidelines and toolkits for the community to effectively scale data for online evaluation within real-world settings in their own scenario. (3) **Cost-effective maintenance to sustain evaluation validity.** Online environment is continuously evolving, making maintaining data validity a challenge. To address this, WebCanvas employs an efficient maintenance strategy with scheduled monitoring and automated alerts that quickly identify action sequences and key nodes validity. When data shifts occur, our test report with error messages guide data owner through quick and effective data corrections. This approach allows us to dynamically adjust our evaluation sets in response to real-time changes in web content with acceptable cost.

Based on WebCanvas framework, we create **Mind2Web-Live** dataset for the community. This dataset contains 542 tasks sampled from Mind2Web (Deng et al., 2024), and we annotate each task with key node verification. Extensive comparisons show that GPT-4 with memory and ReAct (Yao et al., 2023) reasoning achieved the best task success rate of 23.1%. In addition, our online evaluation reveals discrepancies with offline settings, demonstrating that models which perform competitively

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

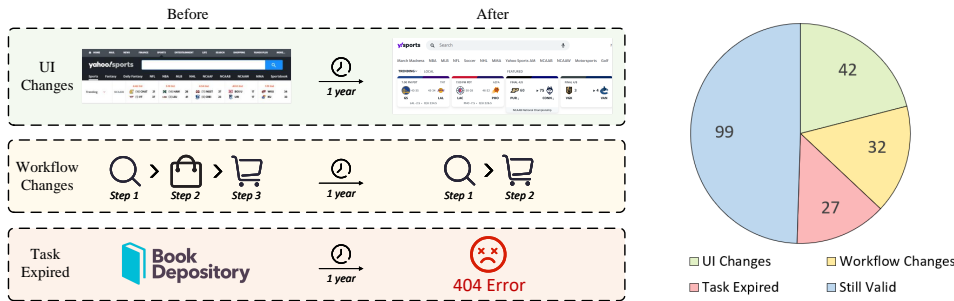


Figure 2: The left side illustrates three typical types of web task evolution over time. We further annotate and classify the status of 200 sampled tasks from Mind2Web training set between May 2023 and May 2024, showing more than half underwent changes to some extent within just one year: 21% experienced UI redesigns, 16% had workflow modifications, and 13.5% became completely expired.

in static offline evaluations do not necessarily maintain their competency in dynamic online environments. We further analyze the impact of various factors specific to online evaluation, such as IP location variability, and suggest maintaining a consistent setup within our framework to ensure reliable results. Finally, we investigate the use of key node annotations as a form of intermediate reward for in-context reasoning. Our findings suggest that web agents can benefit from human-provided reward signal, whereas even advanced models exhibit inaccuracies when generating such intermediate progress indicators without any reference. These inaccuracies subsequently impair execution performance.

2 PROBLEM FORMULATION OF INTERACTIVE WEB-BASED TASK

The real-world web environment can be formulated as: $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{O})$ with state space \mathcal{S} , action space \mathcal{A} (Table 10), deterministic transition function $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ and a state observation space \mathcal{O} (§5). Given a task instruction i , current observation $o_t \in \mathcal{O}$ and the action history $a_{1:t-1}$, an agent issues an action $a_t \in \mathcal{A}$. Consequently, after the execution of the action, the environment transitions to a new state $s_{t+1} \in \mathcal{S}$, and the corresponding observation updates to $o_{t+1} \in \mathcal{O}$. To measure the completion of tasks, we have defined key nodes and evaluation metrics, which are elaborated in §3.1 and §3.2.

3 WEBCANVAS: AN ONLINE EVALUATION FRAMEWORK FOR WEB AGENTS

3.1 DEFINITION OF KEY NODES

The concept of “key nodes” is one of the pivotal ideas in our work. Key nodes refer to indispensable steps in the process of completing specific web tasks, meaning that regardless of the path taken to accomplish a task, these steps are required. These may involve navigation to certain webpages or the performance of specific actions on web pages, such as filling out forms or clicking buttons. This design philosophy not only reflects the dynamic nature of the web environment but also captures the diversity of paths present in real-world web pages.

As illustrated in Figure 1, consider the task of “Find top-rated upcoming adventure movies on Rotten Tomatoes” as an example. Users might start directly at the Rotten Tomatoes homepage or use a search engine to navigate straight to the “New Movies Coming Soon” page of the Rotten Tomatoes. Moreover, when filtering the movies, users might choose to first apply a filter for the “adventure” genre and then sort by popularity, or alternatively, sort by popularity before applying the genre filter. Despite the availability of different paths to achieve the goal, entering the specific page and performing the genre and popularity sorting are essential steps in accomplishing the task. Therefore, these three steps are identified as “key nodes”.

In the dynamic and noisy real-world web environment, identifying these key nodes is challenging due to the potential changes in page content and UI updates, which could render element selector paths obsolete. Therefore, we preferred to use URL state as identifiers for key nodes rather than element interaction, which enhanced the Benchmark’s robustness against layout changes. Only element class

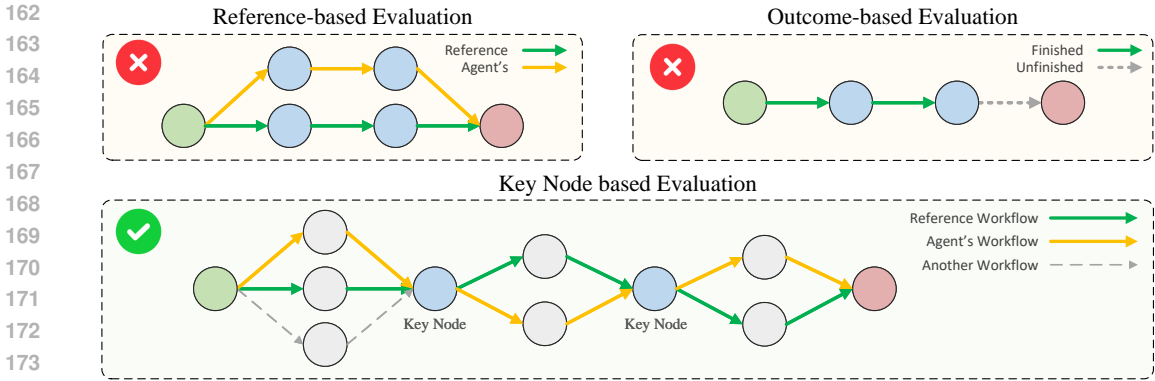


Figure 3: Comparison of different evaluation methods.

methods are considered for key nodes that cannot be represented by URLs. The detailed judgment method is described in Appendix C. By defining key nodes, WebCanvas is able to dynamically assess the execution capabilities of web agents in real-world web environments, offering a practical and flexible evaluation method for the development of web agents.

3.2 EVALUATION METRICS

The evaluation metrics of WebCanvas comprised of two main components: step score and task score. The step score evaluates the agent’s performance with regard to each key node, defining three types of evaluation targets along with three evaluation functions at each step. The task score includes two functions to assess the task’s completeness and overall execution efficiency.

Step Score Inspired by previous works (Zhou et al., 2023; Koh et al., 2024), we introduced three evaluation targets in calculating step score, allowing us to examine from different aspects: **URL**, **Element Path**, and **Element Value**. We implemented three match functions for these targets: **Exact Match**, **Include Match**, and **Semantic Match**. Each key node is associated with an evaluation function, which comprises an evaluation target and a match function. One step score is awarded when the agent successfully reaches a key node and passes the associated evaluation function verification. Table 1 shows a list of applicable evaluation functions and their introductions for reference. Some examples are shown in Table 7. We use **Completion Rate** to represent the MICRO average of Step Scores.

Table 1: Overview of evaluation functions. “E” is short for Web Element.

Match Function	Description	URL	E. Path	E. Value
Exact Match	Precise matching, such as URL parameters or form fields.	✓	✓	✓
Include Match	Evaluates if output includes the reference, ideal for keyword detection.	✓	✗	✓
Semantic Match	Uses LLM for complex content reasoning tasks, like product identification.	✓	✗	✓

Task Score Task performance is evaluated using two main metrics: the **Task Success Rate** and the **Efficiency Score**. The **Task Success Rate** is determined by the proportion of tasks in the test set that are completed successfully, where a task is considered complete if all designated key nodes are achieved. The **Efficiency Score** is calculated using the formula $ES = \frac{L}{P}$. L represents the total number of steps the agent took to complete the task, and P is the cumulative step score obtained by the agent upon completing the task.

4 MIND2WEB-LIVE: A REAL-TIME ONLINE BENCHMARK FOR WEB AGENTS

4.1 DATASET CONSTRUCTION

To develop a real-world online benchmark for web agents, we introduce Mind2Web-Live, which is derived from tasks present in the Mind2Web (Deng et al., 2024) dataset. We employed WebCanvas framework as a guidance for the sampling and re-annotation of these tasks. Consequently, we selectively excluded all tasks that contained time-sensitive descriptions, such as those involving specific dates or times. We randomly sampled 601 tasks from the training set and included all 179 tasks from the cross-task test set, which were then re-annotated in a real-world online environment.

The annotation process presented multiple challenges. Notably, due to updates in website content and operational changes, we discovered 96 tasks that were no longer applicable and subsequently removed them from the dataset. Additionally, 142 tasks were discarded due to ambiguous task definitions, log-in requirements or the difficulty in clearly defining key nodes. To enhance the clarity and reliability of task execution, we revised the instructions for 51 tasks.

For human trajectory and key node annotation, we developed and made public a browser plugin and an annotation platform. After a rigorous annotation and review process, described in Appendix A, 542 high-quality tasks were established for the Mind2Web-Live dataset, including 438 of the training set and 104 of the test set. As shown in Table 2, Mind2Web-Live encompasses 2439 key nodes and 4550 detailed annotation steps. The tasks in the dataset cover a wide range of webpage types and operations, designed to comprehensively evaluate the performance of web agents in a dynamic and variable online environment. The distribution of the evaluation function is illustrated in Figure 5.

4.2 DATASET MAINTENANCE

We pay special attention to the dynamic nature of the benchmark to adapt to the constantly changing web environment. We recognize that updates and changes to website content, such as UI updates, database changes, or website close-down, are inevitable as time progresses. Such changes may lead to the obsolescence of previously defined tasks or key nodes.

We are thus committed to process a regular data maintenance schedule every two months, as shown in Figure 4. We first developed a community-driven platform where dataset users can visualize details of each task and report any issues through a bug-reporting feature. In addition to community supervision, we leveraged the data stored during the annotation stage to ensure a stable playback of these recorded human trajectories, with any invalidities in the workflows or key nodes being promptly reported. Appendix H provides an example of a report. Suspicious tasks are re-annotated during our reviews to ensure that each task accurately reflects the current web environment.

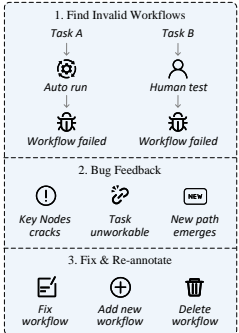


Figure 4: System of maintenance

Over the past two months, we reviewed 104 tasks in the Mind2Web-Live test set. During this review, we found that 5 tasks had underwent key nodes degradation. The four authors took responsibility for the maintenance work, with each spending less than an hour per maintenance cycle, making this an acceptable cost. For invalid workflows, we updated both the trajectories and the key nodes, a process that usually takes about 5 minutes due to our previous annotation experience. For invalid key nodes, we only needed to update the key node functions with around 2 minutes.

5 EXPERIMENT

Inspired by previous work (Yao et al., 2023; Zhou et al., 2023; Zheng et al., 2024), we built a universal baseline agent framework that consists four key modules: Planning, Observation, Memory and Reward. This framework is engineered to be plug and play, operating in real-world online web environments, serving as a foundation for the community to benchmark with rather than introducing new innovations. Detailed implementation is provided in the Appendix E.

Table 2: Data distribution

Statistic	Number
Total selected tasks	780
- Expired Tasks	96
- Unable to annotate	142
- Mind2Web-Live	542
- training set	438
- test set	104
Annotate steps	4550
Avg. steps	8.39 / task
Eval functions	2439
Avg. Eval functions	4.5 / task

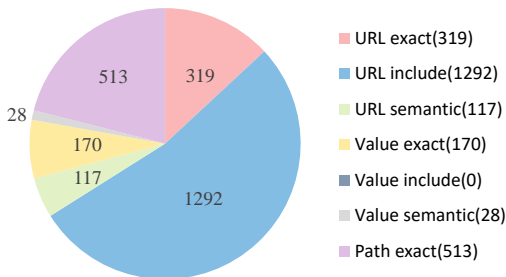


Figure 5: Evaluation Function distribution

Table 3: Comparison of web agent performance in online and offline evaluations. We randomly sampled 40 instances from the Mind2Web-Live test set. These were then tested in both online and offline settings. GPT-3.5 denotes gpt-3.5-turbo-0125, and GPT-4 denotes gpt-4-0125-preview. ‘Task SR(0)’ and ‘Task SR(1)’ denote the Task Success Rates with zero tolerance and tolerance for error at one step (or key node), respectively.

Model	Offline			Online		
	Step SR(%)	Task SR(0)(%)	Task SR(1)(%)	Completion Rate(%)	Task SR(0)(%)	Task SR(1)(%)
MindAct	44.3	10.0	25.0	25.5	7.50	12.5
GPT-3.5	15.5	2.50	7.50	35.4	10.0	17.5
GPT-4	28.4	5.00	22.5	41.1	10.0	25.0

5.1 DISCREPANCY BETWEEN OFFLINE EVALUATION AND ONLINE GENERALIZATION

The settings of evaluation on offline datasets that reflect real-world intents, such as Mind2Web (Deng et al., 2024), are inherently different from WebCanvas framework. Nevertheless, we managed to study the qualitative discrepancy between offline evaluation and online generalization. During online inference, we attempted to reproduce the original prompt of the MindAct model, which was trained and evaluated on the offline dataset, as proposed in the Mind2Web paper. It is important to note that the evaluation metrics used in offline evaluation differ from those proposed in our online evaluation framework. The Step Success Rate in offline testing assesses the accuracy of single-step action prediction, and for the entire task dimension, a positive reward is given only when all single-step actions are correctly predicted, which is not the case in online evaluation, as we evaluate the intermediate state, not the referenced action.

5.2 MAIN RESULT

In our experiments, we observed significant discrepancies between offline and online testing. The results, as detailed in Table 3, show that the model trained on the Mind2Web training set struggles to generalize to the online environment a year later. The comparative performance of MindAct-Large (Deng et al., 2024), GPT-3.5, and GPT-4 in the online environment was opposite to that in offline testing. We further analyzed why such differences occur. Specifically, through human annotation of 100 evaluation steps of GPT-4, we noted that approximately 30% of its actions were reasonable, indicating that these models are possibly capable of generating valid paths under the prompt, albeit penalized by the offline reference-based evaluation method. We also discovered that the MindAct model, when reasoning in an online environment, frequently encountered difficulties in recovering from erroneous states. When entering web pages less relevant to the task goal, the MindAct model had a high probability of outputting null actions, causing the task to terminate.

Furthermore, comparative performance of different models in Table 4 indicates that GPT-4 outperforms other models in both effectiveness and efficiency in web agent tasks within a live environment, with Qwen being the best-performing open-source model. However, there remains considerable room for future enhancements across all models. These results underscore the need for models that can better generalize to dynamic, real-world web environments.

Table 4: Performance of different models without the reward module on the Mind2Web-Live test set, sorted by Completion Rate from highest to lowest. Qualitative analysis of agent performance in online environment are illustrated in Appendix G.

Model	Completion Rate (%)	Task SR (%)	Efficiency Score
GPT-4-0125-preview	48.8	23.1	2.47
Claude-3-Sonnet-20240620	<u>47.9</u>	<u>22.1</u>	2.92
GPT-4o-2024-05-13	47.6	<u>22.1</u>	2.88
Gemini-1.5-Pro	44.6	<u>22.1</u>	4.48
GPT-4-turbo-2024-04-09	44.3	21.2	<u>2.78</u>
Claude-3-Sonnet-20240229	43.9	20.2	3.34
Qwen1.5-110B-Chat	43.9	20.2	4.02
GPT-4o-mini-2024-07-18	42.9	21.2	2.97
DeepSeek-V2	41.2	18.3	4.44
Qwen2-72B-Instruct	40.9	15.4	4.60
Claude-3-Opus-20240229	40.3	14.4	3.52
GPT-3.5-turbo-0125	40.2	16.3	3.03
Mixtral-8x22B	37.2	17.3	4.80
Qwen1.5-72B-Chat	35.6	15.4	4.29
Gemini-Pro	35.3	13.5	4.69
Claude-3-Haiku-20240307	33.4	16.3	4.27
Qwen1.5-7B-Chat	24.5	10.6	8.34

6 ANALYSIS

6.1 FACTORS INFLUENCING AGENT PERFORMANCE

In this section, we delve into the factors influencing agent performance across a range of web tasks. Through a series of experiments, we assessed the impact of task complexity, website dynamics, task domain, key node distribution in the dataset, and the experimental setup—including system specifications, browser engine, and IP location.

Our findings reveal that increased task complexity directly correlates with diminished agent performance. The domain of the task also significantly affects performance, with agents handling entertainment-related tasks more adeptly than those involving shopping or travel. This variation suggests LLMs’ capacity of semantic understanding and reasoning differs across domains and websites. Moreover, the experimental environment plays a crucial role in agent performance. We recommend experimenting on a Windows platform using Chrome or Firefox browser engines, preferably on servers located in the United States. Statistics and experiment results are detailed in Appendix F.2.

6.2 NECESSITY OF KEY NODE EVALUATION IN LIVE ENVIRONMENTS

Previous agent evaluation methods primarily focus on two aspects: reference-based evaluation (Deng et al., 2024; Zheng et al., 2024) and outcome-based evaluation (Zhou et al., 2023; Koh et al., 2024; Mialon et al., 2023). However, these methods falter when applied to the unpredictable nature of live web tasks. To address the inherent variability in task completion paths within an online evaluation framework, we employed Sankey diagrams to visualize the trajectories of our web agent and human demonstrations on tasks where our agent successfully navigated all designated key nodes in Figure 10 within §F.2.

We further annotate Mind2Web-Live test set to identify whether the final key node is a sufficient condition for task completion. It turns out only 46 out of 104 tasks met this criterion. This finding starkly illustrates that solely evaluating the final state or outcome is inadequate for web environments that are not fully reproducible. As shown in Figure 3, key node based evaluation enhances explainability of agent performance, prevents illegal shortcuts taken and facilitates the modeling of structured in-progress rewards, valuable for both in-context reasoning experiments and future reinforcement learning training.

Table 5: Performance of different models with reward module, based on a random sample of 130 tasks from the Mind2Web-Live dataset. “(+)” indicates the inclusion of a reward module with human-labeled reward. Bold numbers represent the best values across different planning models. Model notation follows Table 3, except for gpt-4-vision-preview(GPT-4V). Human Alignment score represents agents’ alignment with human decision on task completion, while the larger indicates better alignment, detailed in Appendix D.

Planning Model	Reward Model	Completion Rate (%)	Task Success Rate (%)	Efficiency Score	Human Alignment
GPT-3.5	/	34.6	13.8	5.25	/
GPT-4	/	46.9	16.9	3.77	/
GPT-4	GPT-3.5	43.5	16.2	3.24	0.445
GPT-4	GPT-4	42.1	13.8	3.07	0.430
GPT-3.5	GPT-4	36.6	10.8	3.73	0.385
GPT-4	GPT-4V	42.4	8.5	3.42	0.419
GPT-3.5	GPT-4(+)	43.6	13.8	3.28	0.452
GPT-4	GPT-4(+)	52.3	12.3	3.27	0.506
GPT-4	GPT-4V(+)	51.3	12.3	2.71	0.502

6.3 PLANNING WITH HUMAN-LABELED REWARD

Reward modeling for agent tasks is crucial in both in-context reasoning and reinforcement learning (Shinn et al., 2024; Bai et al., 2024). Previous research has proven that LLMs can generate high quality reward signal to enhance reasoning performance across various agent tasks (Shinn et al., 2024; Pan et al., 2024). However, recent research adopting an un-tuned foundation model for self-reward prediction shows that their effectiveness is not consistent in specific domains (Olausson et al., 2023; Shinn et al., 2024). Our preliminary experiments indicate that agent performance do not benefit from a self-reward module in the online web environment. This is attributed to several factors, such as overconfidence in task completion assessments and the long-term impact of poor-quality reward signals accumulated in agent memory. Thus it raises a natural question - *Does the quality of the reward signal hinder the self-reward module’s effectiveness in online web environments?* In our study, we introduced a reward module with human-labeled reward. The experimental results on Mind2Web-Live, which confirm our hypothesis, are detailed in Table 5.

From the original data, we extracted post-action URLs, action types, CSS selector paths, and key nodes functions as metadata for our golden reference synthesis. We then employed a carefully designed prompt available in Appendix K, using GPT-4 to generate a structured linguistic guidance for task progress estimation for each task. This guidance includes the overall goal of the task and task completion criteria, specifically highlighting all key nodes for the task to be considered fully completed. We then integrate the content of the current task’s golden reference with the original design of history and current observation for reward reasoning. From comprehensive experiments, we find that the integration of a reward module does not enhance agent performance and may even lead to a decline in Task Success Rate and Task Completion Rate. This finding aligns with findings in (Shinn et al., 2024) about the effect of self-reflection modules in web agent tasks. However, we find the Completion Rate improves in both GPT-3.5 and GPT-4 experiments with the integration of a reward module with human-labeled reward, despite the reward module triggering premature stops. These findings point out the importance of better reward modeling in web agent reasoning.

7 RELATED WORKS

Agent Benchmarks Early researches (Shi et al., 2017) (Liu et al., 2018) provided relatively simple simulations and assessment methods for web navigation tasks. However, with the rise of Large Language Models, these methods have become inadequate for assessing agents’ capability. Recent studies have chosen to construct realistic simulated environments (Yao et al., 2022; Zhou et al., 2023; Koh et al., 2024; Drouin et al., 2024), use offline saved datasets (Deng et al., 2024; Lu et al.,

Table 6: Case study of previous benchmarks

Benchmark	Real-world Intents	Dynamic Environment	Keep Updated	Intermediate Env. State	Easy to Scale	Disk Usage
MiniWoB++ (Liu et al., 2018)	✗	✓	✓	✗	✗	< 1GB
WebShop (Yao et al., 2022)	✗	✓	✗	✗	✗	~ 10GB
Mind2Web (Deng et al., 2024)	✓	✗	✗	✗	✗	~ 10GB
WebArena (Zhou et al., 2023)	✓	✓	✗	✗	✓	> 100GB
VWebArena (Koh et al., 2024)	✓	✓	✗	✗	✓	> 100GB
GAIA (Mialon et al., 2023)	✓	✓	✗	✗	✓	< 1GB
WEBLINX (Lu et al., 2024)	✓	✗	✗	✓	✗	< 1GB
OmniACT (Kapoor et al., 2024)	✓	✗	✗	✓	✗	< 1GB
WebCanvas	✓	✓	✓	✓	✓	< 1GB

2024), or select relatively stable answers to assess the capabilities of web agents (Mialon et al., 2023). In terms of dynamic evaluation methods, many studies (Kiela et al., 2021; Ma et al., 2021; Jain et al., 2024) have proposed their own solutions. Moreover, beyond network platforms, several initiatives have also been undertaken on other platforms such as Android mobile devices, operating systems, and databases (Rawles et al., 2024; Liu et al., 2024b; Xie et al., 2024). We perform a more comprehensive case study on previous web agent benchmarks in Table 6, WebCanvas aims to more comprehensively test agents’ capability in the real world through key nodes and corresponding evaluation functions.

Agent Frameworks In the area of reasoning frameworks, several studies have achieved notable success in logical reasoning challenges (Wei et al., 2022; Yao et al., 2024; 2023; Shinn et al., 2024; Summers et al., 2024). Regarding web agent reasoning frameworks, many researches has been conducted to enhance the capabilities of web agents (Nakano et al., 2021; Gur et al., 2024; Gür et al., 2023; Kim et al., 2024; Lo et al., 2023; Lai et al., 2024). Some studies have introduced multi-modal modules that integrate visual and semantic information, thereby enhancing the capabilities of agents on web platforms (Zheng et al., 2024; Furuta et al., 2024; He et al., 2024).

8 DISCUSSION & LIMITATIONS

Developing a suitable evaluation framework is a fundamental component in the advancement of autonomous web agents. This research addresses the challenge of live evaluation in a real-world web environment. Among these are the need to define key nodes in a completely open environment, unify the inference processes across different digital autonomous agents, and reduce the maintenance costs associated with real-time data and evaluation functions. Through our efforts, we have made significant strides toward establishing a robust and accurate online evaluation system for web agents.

However, the transition to live, dynamic evaluations in unpredictable online environments introduces new complexities not present in controlled, offline settings. The unsolved challenges we encountered in online evaluation of web agents include network instability, dynamic and complex task pathways, and the limitations of static evaluation functions. These challenges highlight the necessity for ongoing research and community efforts to refine and enhance evaluation frameworks for autonomous web agents in complex, real-world environments. For more details, please refer to Appendix J.

9 CONCLUSION

In this work, we have pioneered the development of framework for online evaluation of web agents, and investigated the challenges associated with online evaluation and the difficulties faced by current web agent reasoning frameworks in online inference. Simultaneously, we have constructed toolkits and a community-driven platform that empowers web agent researchers and developers to build datasets and evaluate their web agent frameworks and models in an online environment while collecting feedback on dataset design, data annotation quality, and data validity throughout the process. We strongly encourage further work on online datasets, web agents, and evaluation function designs. By fostering a collaborative and iterative value to dataset creation and evaluation, we eagerly anticipate the continued growth of advancement of autonomous intelligence.

486 ETHICS STATEMENT

487

488 In this work, all annotators for the Mind2Web-Live dataset were fully informed about the data
 489 annotation tasks and compensated for their participation. We take full care of privacy and data
 490 security, ensuring that no sensitive personal information was collected, and all tasks can be completed
 491 without login.

492

493 REPRODUCIBILITY STATEMENT

494

495 We have taken significant measures to ensure the reproducibility of our results. The code used for the
 496 experiments, as well as the Mind2Web-Live dataset, are available in the supplementary materials.
 497 Additionally, comprehensive guidelines for setting up the environment and the parameters required to
 498 run the experiments are included.

499

500 REFERENCES

501

502 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
 503 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
 504 *arXiv preprint arXiv:2303.08774*, 2023.

505

506 Hao Bai, Yifei Zhou, Mert Cemri, Jiayi Pan, Alane Suhr, Sergey Levine, and Aviral Kumar. Digirl:
 507 Training in-the-wild device-control agents with autonomous reinforcement learning. *arXiv preprint*
 508 *arXiv:2406.11896*, 2024.

509

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

513 Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su.
 514 Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing*
 515 *Systems*, 36, 2024.

516

517 Alexandre Drouin, Maxime Gasse, Massimo Caccia, Issam H Laradji, Manuel Del Verme, Tom
 518 Marty, David Vazquez, Nicolas Chapados, and Alexandre Lacoste. Workarena: How capable are
 519 web agents at solving common knowledge work tasks? In *Forty-first International Conference on*
 520 *Machine Learning*, 2024.

521

522 Hiroki Furuta, Kuang-Huei Lee, Ofir Nachum, Yutaka Matsuo, Aleksandra Faust, Shixiang Shane
 523 Gu, and Izzeddin Gur. Multimodal web navigation with instruction-finetuned foundation models.
 524 In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=efFmBWioSc>.

525

526 Izzeddin Gür, Ofir Nachum, Yingjie Miao, Mustafa Safdari, Austin Huang, Aakanksha Chowdhery,
 527 Sharan Narang, Noah Fiedel, and Aleksandra Faust. Understanding html with large language
 528 models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 2803–
 529 2821, 2023.

530

531 Izzeddin Gur, Hiroki Furuta, Austin V Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and
 532 Aleksandra Faust. A real-world webagent with planning, long context understanding, and program
 533 synthesis. In *The Twelfth International Conference on Learning Representations*, 2024. URL
 534 <https://openreview.net/forum?id=9JQttrumvg8>.

535 Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan,
 536 and Dong Yu. WebVoyager: Building an end-to-end web agent with large multimodal models. In
 537 *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume*
 538 *1: Long Papers)*, pp. 6864–6890, Bangkok, Thailand, August 2024. Association for Computational
 539 Linguistics. doi: 10.18653/v1/2024.acl-long.371. URL <https://aclanthology.org/2024.acl-long.371>.

- 540 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando
541 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free
542 evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- 543
544 Raghav Kapoor, Yash Parag Butala, Melisa Russak, Jing Yu Koh, Kiran Kamble, Waseem Alshikh,
545 and Ruslan Salakhutdinov. Omniact: A dataset and benchmark for enabling multimodal generalist
546 autonomous agents for desktop and web. *arXiv preprint arXiv:2402.17553*, 2024.
- 547
548 Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie
549 Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, et al. Dynabench: Rethinking bench-
550 marking in nlp. In *Proceedings of the 2021 Conference of the North American Chapter of the*
551 *Association for Computational Linguistics: Human Language Technologies*, pp. 4110–4124, 2021.
- 552
553 Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language models can solve computer tasks.
554 *Advances in Neural Information Processing Systems*, 36, 2024.
- 555
556 Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Lim, Po-Yu Huang, Graham Neubig,
557 Shuyan Zhou, Russ Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating multimodal
558 agents on realistic visual web tasks. In *Proceedings of the 62nd Annual Meeting of the Association*
559 *for Computational Linguistics (Volume 1: Long Papers)*, pp. 881–905, Bangkok, Thailand, August
2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.50. URL
<https://aclanthology.org/2024.acl-long.50>.
- 560
561 Hanyu Lai, Xiao Liu, Iat Long Iong, Shuntian Yao, Yuxuan Chen, Pengbo Shen, Hao Yu, Hanchen
562 Zhang, Xiaohan Zhang, Yuxiao Dong, et al. Autowebglm: Bootstrap and reinforce a large language
563 model-based web navigating agent. *arXiv preprint arXiv:2404.03648*, 2024.
- 564
565 Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, and Percy Liang. Reinforcement learning on
566 web interfaces using workflow-guided exploration. In *International Conference on Learning*
Representations (ICLR), 2018.
- 567
568 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in*
neural information processing systems, 36, 2024a.
- 569
570 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,
571 Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui
572 Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang.
573 Agentbench: Evaluating LLMs as agents. In *The Twelfth International Conference on Learning*
Representations, 2024b. URL <https://openreview.net/forum?id=zAdUB0aCTQ>.
- 574
575 Robert Lo, Abishek Sridhar, Frank F Xu, Hao Zhu, and Shuyan Zhou. Hierarchical prompting
576 assists large language model on web navigation. In *Findings of the Association for Computational*
Linguistics: EMNLP 2023, pp. 10217–10244, 2023.
- 577
578 Xing Han Lu, Zdeněk Kasner, and Siva Reddy. WebLINX: Real-world website navigation with
579 multi-turn dialogue. In *Forty-first International Conference on Machine Learning*, 2024. URL
580 <https://openreview.net/forum?id=mUSPhG4uDW>.
- 581
582 Zhiyi Ma, Kawin Ethayarajh, Tristan Thrush, Somya Jain, Ledell Wu, Robin Jia, Christopher
583 Potts, Adina Williams, and Douwe Kiela. Dynaboard: An evaluation-as-a-service platform for
584 holistic next-generation benchmarking. *Advances in Neural Information Processing Systems*, 34:
585 10351–10367, 2021.
- 586
587 Grégoire Mialon, Clémentine Fourier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom.
588 Gaia: a benchmark for general ai assistants. *arXiv preprint arXiv:2311.12983*, 2023.
- 589
590 Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher
591 Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted
592 question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- 593
Theo X Olausson, Jeevana Priya Inala, Chenglong Wang, Jianfeng Gao, and Armando Solar-Lezama.
Is self-repair a silver bullet for code generation? In *The Twelfth International Conference on*
Learning Representations, 2023.

- 594 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
595 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
596 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–
597 27744, 2022.
- 598 Jiayi Pan, Yichi Zhang, Nicholas Tomlin, Yifei Zhou, Sergey Levine, and Alane Suhr. Autonomous
599 evaluation and refinement of digital agents. In *First Conference on Language Modeling*, 2024.
600 URL <https://openreview.net/forum?id=NPAQ6FKSmK>.
- 601 Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. An-
602 droidinthewild: A large-scale dataset for android device control. *Advances in Neural Information*
603 *Processing Systems*, 36, 2024.
- 604 Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Percy Liang. World of bits: An
605 open-domain platform for web-based agents. In *International Conference on Machine Learning*,
606 pp. 3135–3144. PMLR, 2017.
- 607 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion:
608 Language agents with verbal reinforcement learning. *Advances in Neural Information Processing*
609 *Systems*, 36, 2024.
- 610 Theodore Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas Griffiths. Cognitive architectures
611 for language agents. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL
612 <https://openreview.net/forum?id=li6ZCvf1QJ>. Survey Certification.
- 613 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
614 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
615 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 616 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
617 Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
618 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 619 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
620 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
621 *neural information processing systems*, 35:24824–24837, 2022.
- 622 Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing
623 Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal
624 agents for open-ended tasks in real computer environments. *arXiv preprint arXiv:2404.07972*,
625 2024.
- 626 S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, , and Y. Cao. React: Synergizing reasoning
627 and acting in language models. In *International Conference on Learning Representations (ICLR)*,
628 2023.
- 629 Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable
630 real-world web interaction with grounded language agents. *Advances in Neural Information*
631 *Processing Systems*, 35:20744–20757, 2022.
- 632 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.
633 Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural*
634 *Information Processing Systems*, 36, 2024.
- 635 Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. GPT-4v(ision) is a generalist web
636 agent, if grounded. In *Forty-first International Conference on Machine Learning*, 2024. URL
637 <https://openreview.net/forum?id=piecKJ2D1B>.
- 638 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,
639 Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for building
640 autonomous agents. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*, 2023.
641
642
643
644
645
646
647

A DATA COLLECTION DETAILS

A.1 RECORDING PROCESS

In the construction of Mind2Web-Live, the quality and reliability of the data are paramount. To this end, we have employed an efficient tool (Figure 6) for recording browser operations. The tool precisely captures browser interaction from the users, covering a wide range of activities such as clicks and input actions. The recorded details include the type of operation, execution parameters, target element’s selector path, element content, and its coordinates on the webpage. Moreover, the tool accompany each step with a webpage screenshot, not only facilitating process replication but also providing a visual reference for workflow validation and review (Figure 18, 19). This approach enables us to comprehensively record all the steps required to complete specific tasks, forming the foundation of Mind2Web-Live. Upon completion of the data recording, we meticulously annotated the key nodes of each process along with their corresponding Evaluation Functions.

A.2 ANNOTATION PROCESS

In our study, the annotation process plays a pivotal role in ensuring data quality and task validity. To ensure the accuracy and consistency of data annotations, we assembled an annotation team comprised of several authors of this paper and five senior undergraduate students majoring in Computer Science. Not only do the members of the annotation team possess a solid background in Computer Science, but they also received specialized training to ensure consistency in their understanding and identification abilities in annotating key nodes. Prior to beginning the formal annotation process, all annotators were rigorously trained over a period of two weeks, which included trial annotations that were subsequently not included in the final dataset.

During the annotation phase, we employed a comprehensive reward mechanism. Each annotator was compensated based on the number of tasks they completed, with additional bonuses awarded for high-quality annotations to encourage precise and consistent results. This combined reward system not only bolstered work enthusiasm but also enhanced the overall quality of the annotation work, laying a solid foundation for the construction of an efficient web agent benchmark.

To guarantee the quality of annotations, we instituted a variety of strategies. Each task was annotated independently by one annotator, followed by individual reviews by two other members to verify the accuracy of the key nodes. Throughout the annotation process, we regularly organized discussion sessions for the annotation team to share their experiences and challenges encountered, thereby improving the overall efficiency and quality of the team’s annotations.

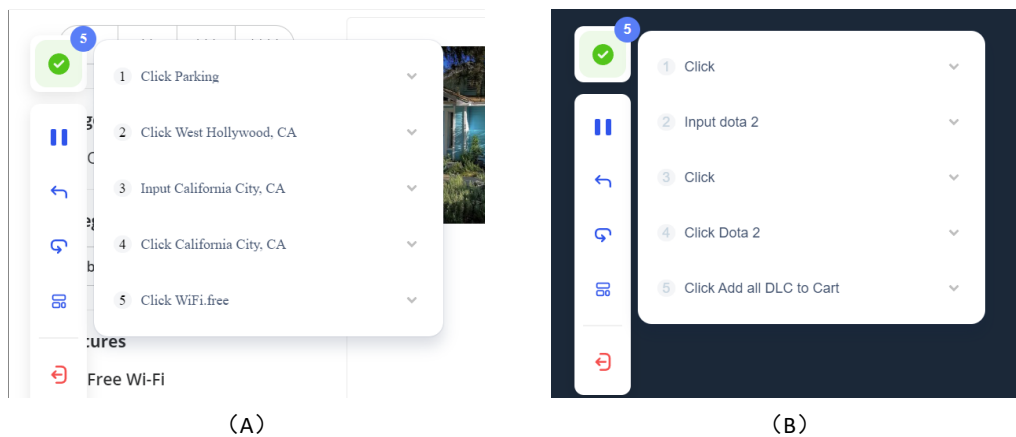
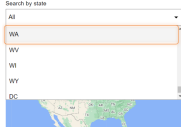
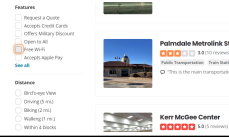



Figure 6: An illustration of the Annotate Tool being used to annotate two distinct tasks: (A) “Find parking in California city for Limos which also offers free wi-fi in yelp”, and (B) “Find Dota 2 game and add all DLC to cart in steam”.

Table 7: Example annotations of the Evaluation Functions

State	Title	Annotation Details
	Locate a large store in Washington that has kids' and maternity products in uniqlo	Evaluation Function: Element value semantic match Instructions: Decide Whether is searching for Washington D.C.
	Find parking in California city for Limos which also offers free wi-fi in yelp	Evaluation Function: URL include match Param: attrrs Value: WiFi.free
	Find Dota 2 game and add all DLC to cart in steam	Evaluation Function: Element path exact match Selector: <code>//*[@id="dlc_purchase_action"] /div[2]/a/span</code>

A.3 TASK DISTRIBUTION AND DOMAIN COVERAGE

See Table 8.

Table 8: Task Distribution and Domain Coverage

Domain	Subdomain	Mind2Web-Live Test	Mind2Web-Live Train
Entertainment	Sports	9	32
	Event	5	20
	Game	3	24
	Movie	9	30
	Music	5	18
Shopping	General	3	28
	Auto	7	33
	Department	6	8
	Digital	6	15
	Fashion	3	15
Travel	Speciality	13	44
	General	0	11
	Airlines	5	18
	Car rental	1	11
	Ground	9	28
Travel	Hotel	3	12
	Restaurant	6	31
	Other	11	60
Total		104	438

B COMPARISON OF THE MIND2WEB-LIVE AND MIND2WEB DATASETS

Table 9: Comparison of the Mind2Web-Live and Mind2Web Datasets. “Ele.” indicates “Element”, “Op.” indicates “Option” and “SR” indicates “success rate”.

Attributes	Mind2Web-Live	Mind2Web
Dataset Size	542	2350
Evaluation Environment	Real-world Online	Offline
Evaluation State	Key Nodes	Each Step
Target Element	Element, URL	Element, Option
Evaluation Metrics	Step Score & Task Score	Step(Ele., Op.) SR & Task SR
Avg. Steps	8.39 / task	7.3 / task

C HOW TO DEFINE EVALUATION FUNCTIONS

For input operations on the page First, determine whether it is a necessary condition for task completion. If it is a necessary condition, then judge whether the execution result can be reflected by the change of the URL. If so, simply take the state after execution as the key node and select the evaluation function as URL exactly/included/semantic match.

If it cannot be reflected by changes in the URL, it needs to be defined as a key node based on click or input operations. Select element path exactly match or element value exactly/included/semantic match for input operations (to determine whether the content of the input element matches).

For click operations on the page Firstly, determine whether it is a necessary condition for completing the task. If it is a necessary condition, then judge whether the execution result can be reflected by the change of the URL. If so, simply take the state after execution as the key node and select the match rule as URL exactly/included/semantic match.

If it cannot be reflected by the change of URL, each click operation should be defined as a key node, and the match can be selected as element element path exactly match or element value match.

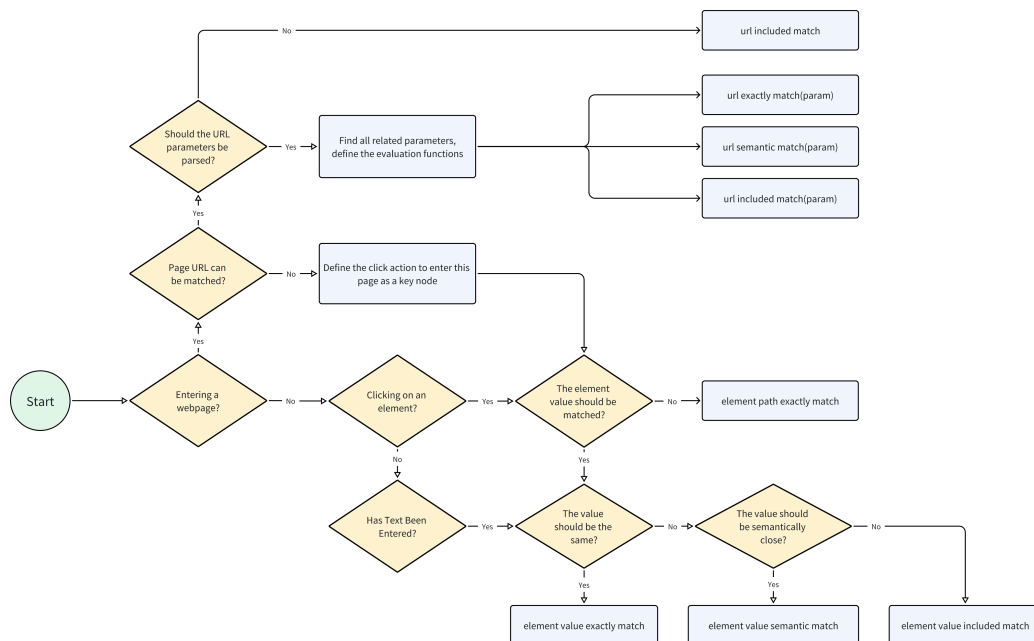


Figure 7: Guidance on how to define an evaluation function for a key node.

D ADDITIONAL EVALUATION METRICS

Human Alignment Score The Human Alignment Score(HAS) assesses how well an agent’s workflow aligns with human behavior. It’s crucial for agents not just to be efficient, but to operate in ways that resemble human actions. The evaluation of this aspect is conducted by contrasting the agent’s task completion signal with the ground truth annotations provided by humans, to gauge the level of consistency. An agent that accurately issues a completion signal upon task completion is deemed to exhibit a high degree of alignment with human behavior, thus earning a full score of one point. Conversely, a delay in issuing the completion signal upon task completion results in a deduction of 0.05 points from the full score as a penalty for decision latency. In instances where an agent stops its operation before accomplishing all the task objectives, the score is determined by the ratio of the step score attained to the maximum step score achievable for that task. Furthermore, if a task is not fully completed and the system forcibly terminates the process due to reaching the maximum step limit, the score awarded is 0.8 times the proportion of the step score attained. The specific algorithm is shown in the formula, where P represents achieved step scores, P_{max} denotes the max step scores of the task.

$$HAS = \begin{cases} 1 & \text{if task is completed with completion signal} \\ 0.95 & \text{if task is completed without completion signal} \\ \frac{P}{P_{max}} & \text{if task is incomplete but completion signal} \\ 0.8 \times \frac{P}{P_{max}} & \text{if task is incomplete and is terminated} \end{cases} \quad (1)$$

E EXPERIMENTAL SETTINGS

E.1 AGENT FRAMEWORK

Planning Integrates past action history, current observations, and task instruction to plan future actions and determine operational values based on the ReAct (Yao et al., 2023) reasoning framework. It can be formally expressed as: $\text{Planning}(\mathbf{h}_{1:t}, \mathbf{o}_t, \mathbf{i}) \rightarrow (\mathbf{z}_t, \mathbf{a}_t)$, where $\mathbf{h}_{1:t}$ represents history information until time t , \mathbf{o}_t is the observation at time t , \mathbf{i} is the task instruction, while the outputs \mathbf{z}_t and \mathbf{a}_t are the thought and action at time t respectively.

Observation Processes the current webpage’s source code and screenshots, producing an accessibility tree (Zhou et al., 2023) and visual observations as \mathbf{o}_t . In our planning model, we solely focus on textual observations, as visual images involve various grounding mechanisms which could detract from the main focus of our paper. We plan to address this aspect in future research.

Memory Responsible for storing the task instruction and tracking the agent’s operational history, including thoughts and actions history across states. It can be formally expressed as $\mathbf{h}_{1:t} = (\mathbf{z}_{1:t}, \mathbf{a}_{1:t}, \mathbf{r}_{1:t})$ within the framework, where $\mathbf{r}_{1:t}$ denotes the history of reward signal if presents.

Reward Utilizes a self-reflection structure (Shinn et al., 2024), providing a series of reward signal, including a verbal reflection and signal on whether the task is completed. This can be formalized as $\text{Reward}(\mathbf{h}_{1:t}, \mathbf{i}, \mathbf{o}_{t+1}) \rightarrow \mathbf{r}_t$.

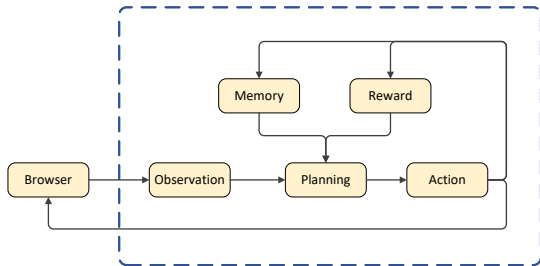


Figure 8: Agent framework

E.2 ACTION SPACE

Table 10: Action space

Action	Operation value
Goto	Value
Google Search	Value
Click	Target id
Hover	Target id
Fill Form	Target id, value
Fill Search	Target id, value
Select	Target id, value
Switch Tab	Target id
Go Back	/

E.3 ADDITIONAL EXPERIMENT SETTINGS

Dataset Sampling Our main experiments were conducted on the Mind2Web-Live test set to avoid data contamination. For experiments involving self-reward, we sampled 130 cases from the complete Mind2Web-Live dataset, ensuring a broad representation free from any dataset-specific biases.

Parameters & Computational Resources The foundation models used across our experiments were standardized with a maximum token of 500 and a temperature setting of 0.7. Computational resources were provided by AWS EC2. While most experiments were conducted on standard compute instances, experiments involving the MindAct model utilized two T4 GPUs to accommodate the model’s computational demands. In addition to using APIs provided by the model developers, our model inference services also incorporated Mixtral-8x22B inference services from Together.ai¹. For the stopping criteria, in experiments with a reward module, we employ reward module to determine whether a process has been completed, otherwise we set a maximum reasoning step length of 1.5 times the annotated task length. Prompts of our experiment can be found in Appendix K.

E.4 OBSERVATION SPACE

Accessibility Tree We employ an accessibility tree-based approach to extract the fundamental textual feature representation from the web environment. The accessibility tree serves as an abstract representation of the structure of a web page, detailing the characteristics of each element within the page. However, the accessibility tree contains a significant amount of redundant information, necessitating the use of a stringent set of filtering criteria to select interactive elements. These filtering criteria include the element’s tag, visibility, usability, as well as textual or image content. Concurrently with the construction of the accessibility tree, we annotate each filtered interactive element, providing information such as element ID, tag, and content. For example, ([1] input ‘search’, etc.). This annotation LLM method facilitates the precise generation of corresponding CSS selector paths during subsequent LLM prediction and execution phases, thereby accurately locating the required elements.

Screenshot We capture screenshots of the current web page to obtain its visual representation and provide this visual context to visual language models, such as GPT-4V. This input method mimics human visual perception, allowing the model to gather the most comprehensive information from the web page. Compared to relying solely on the accessibility tree, using screenshots enhances the ability to identify the layout, appearance, and positioning of web elements more effectively. Additionally, it captures interactive elements and other crucial page information that the accessibility tree might miss. To balance inference costs and recognition effectiveness, the original resolution of the screenshots is set to 1080×720 , though users can define the screenshot resolution according to their specific needs in practical applications.

¹<https://api.together.xyz/models>

F MORE RESULTS OF EXPERIMENTS

F.1 ADDITIONAL MAIN RESULTS

F.1.1 RESULTS ON MIND2WEB-LIVE TRAINING SET

See Table 11.

Table 11: Performance of different models on Mind2Web-Live training set without reward module. As for the model, we experiment with gpt-3.5-turbo-0125 (GPT-3.5), gpt-4-0125-preview (GPT-4).

Model	Completion Rate (%)	Task SR (%)	Efficiency Score
GPT-3.5	<u>34.6</u>	<u>13.8</u>	<u>5.25</u>
GPT-4	46.9	20.1	3.77
Gemini-Pro	31.3	9.23	6.50
DeepSeek-V2	31.8	12.4	5.55
Mixtral-8x22B	29.7	9.44	6.52

F.1.2 ABLATION STUDY

See Table 12.

Table 12: Ablation study on memory and ReAct reasoning architecture (Yao et al., 2023). Results show interesting findings that less capable models like GPT3.5 and Mistral-8x22B do not benefit from memory and advanced reasoning architecture in online web tasks. We encourage more comprehensive evaluation of these modules in web agent framework in future research.

Model	Memory	ReAct	Completion Rate	Task SR	Efficiency Score
GPT-3.5	✓	✓	40.2%	16.5%	3.03
GPT-4	✓	✓	48.8%	23.1%	2.47
Mixtral-8x22B	✓	✓	37.2%	17.3%	4.80
GPT-3.5	✗	✓	43.5%(↑ 3.3%)	19.2%(↑ 2.7%)	3.12(↓ 0.09)
GPT-3.5	✓	✗	42.5%(↑ 2.3%)	22.1%(↑ 5.6%)	2.98(↑ 0.05)
Mixtral-8x22B	✗	✓	42.3%(↑ 5.1%)	17.3%(-)	4.39(↑ 0.41)
Mixtral-8x22B	✓	✗	42.5%(↑ 5.3%)	19.2%(↑ 1.9%)	4.40(↑ 0.40)
GPT4	✗	✓	48.6%(↓ 0.2%)	20.9%(↓ 2.2%)	2.70(↓ 0.23)
GPT4	✓	✗	46.6%(↓ 2.2%)	22.1%(↓ 1.0%)	2.67(↓ 0.20)

F.2 ADDITIONAL ANALYSIS

See Table 13, Figure 9, Figure 10, Figure 11, Figure 12, Figure 13, Figure 14.

Table 13: Experiment on IP Regions and devices. It presents the results of experiments conducted using the GPT-3.5 planning model across different IP regions, systems and devices. We recommend experimenting on a Windows server using Chrome or Firefox browser engines, preferably on servers located in the United States or Singapore.

Planning Model	IP Region	System	Browser	Completion Rate	Task Success Rate	Efficiency Score
GPT-3.5	United States	Windows	Chrome	40.2%	16.5%	3.03
GPT-3.5	United States	Windows	Firefox	42.1%	20.2%	2.79
GPT-3.5	United States	Linux	Chrome	36.5%	15.4%	3.33
GPT-3.5	United Kingdom	Windows	Chrome	23.6%	8.65%	7.78
GPT-3.5	Singapore	Windows	Chrome	42.3%	21.2%	2.95

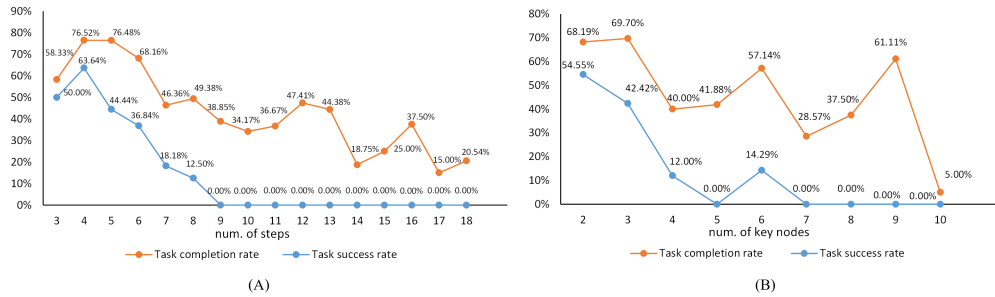
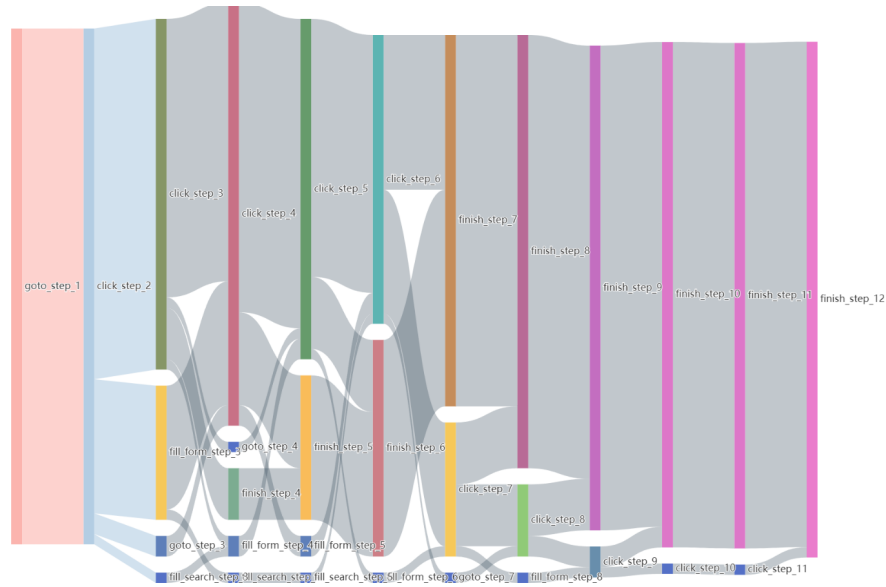


Figure 9: The relationship between task complexity and task difficulty. The “step count” refers to the length of the action sequence in the annotated data, which, along with the number of key nodes, serves as a reference for task complexity.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

Sankey diagram with annotation data



Sankey diagram with Agent's success task data

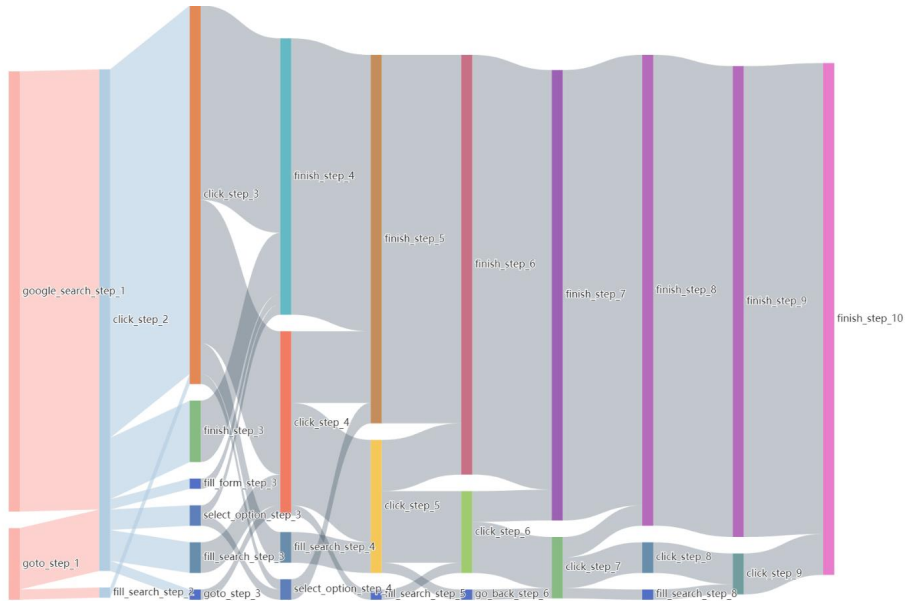


Figure 10: Sankey diagram comparing human demonstration trajectories(A) and agent’s trajectories(B). We randomly sampled 50 success tasks from GPT-4 based agent on the Mind2Web-Live training and testing set to analyze the discrepancy between these trajectories.

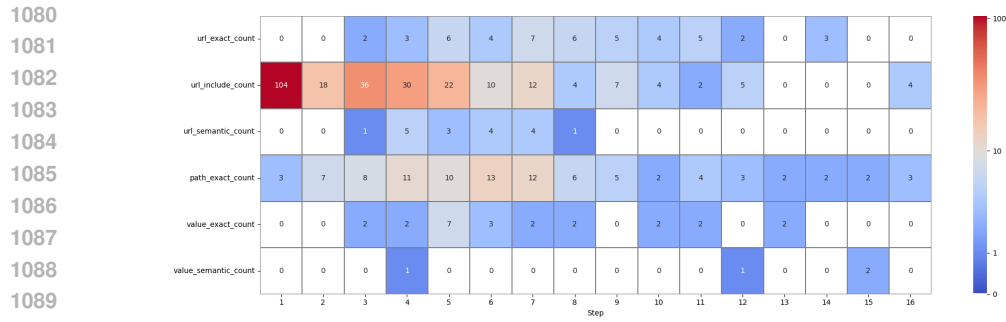


Figure 11: Heatmap of evaluation function counts over annotation steps for the Mind2Web-Live test set. It shows logarithmically transformed counts over various steps. White represents a count of 0, blue indicates smaller counts, and red indicates larger counts. The logarithmic scale helps to evenly distribute the color intensity for better visualization.

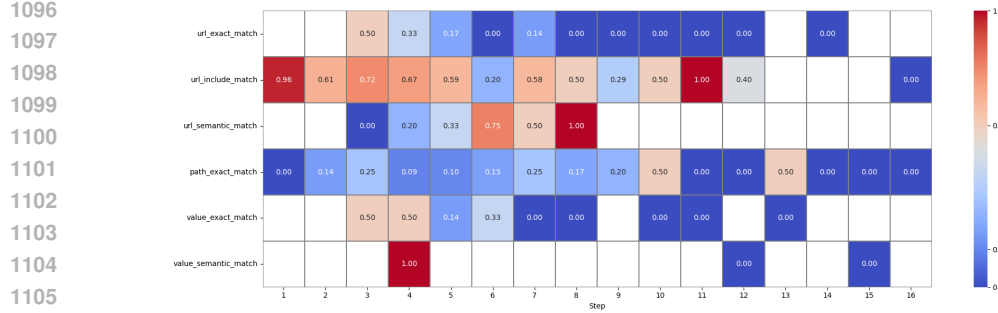


Figure 12: Heatmap of evaluation function accuracy over annotation steps for the Mind2Web-Live test set. The experimental data is derived from GPT-4’s performance on the test sets. The heatmap displays logarithmically transformed accuracy of evaluation functions across different steps. Blue indicates lower accuracy, while red indicates higher accuracy.

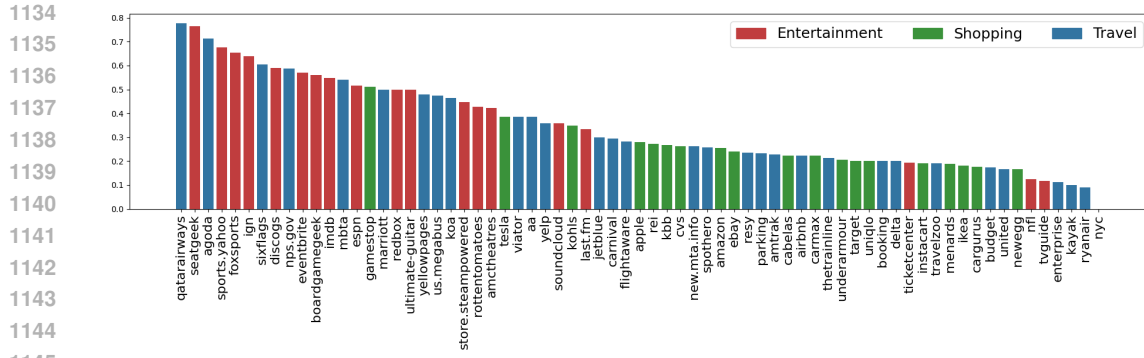


Figure 13: Completion Rate of different website tasks. Due to the large number of websites and the limited number of tasks in the test set, the experimental data is derived from GPT-4’s performance on both the training and test sets. We encourage the community to collaborate in gathering data on online web agent execution across specific websites and tasks.

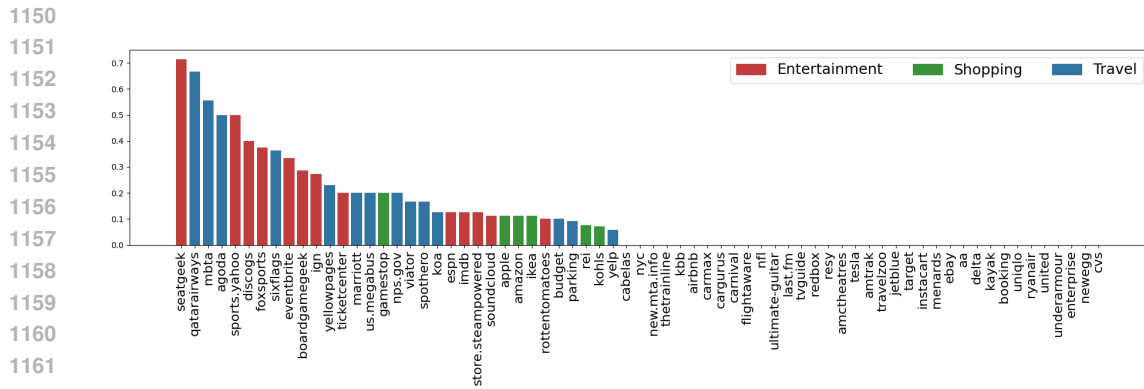


Figure 14: Task Success Rate of different website tasks. Due to the large number of websites and the limited number of tasks in the test set, the experimental data is derived from GPT-4’s performance on both the training and test sets.

1188 G QUALITATIVE ANALYSIS OF EXPERIMENTS

1189

1190 In this section, we conducted a qualitative analysis of error cases in our experimental results. Typical
1191 errors include: local optima, premature termination of tasks, and information loss during inference.
1192

1193

1193 G.1 LOCAL OPTIMA

1194

1195 In our online environment experiments, a task may involve multiple constraints or requirements.
1196 Web pages often contain numerous clickable links, and frequently feature interactable elements
1197 with similar or even identical names. Due to a lack of prior knowledge about the web domain
1198 associated with current task and confusion caused by similar elements, the planning module’s local
1199 decision-making for the current web state is not always accurate. Moreover, our web agent lacks
1200 proactive thinking to revert to an intermediate state within a limited number of steps, thus stuck in a
1201 local optima of the task. This is one of the main reasons for the low task success rate. As shown in
1202 the first line in Table 14, in the task “Check the rating and user reviews for the game ‘Deathloop’ on
1203 IGN”, the web agent ended up at the review article page for ‘Deathloop’ on IGN due to incorrect
1204 path selection from the Google search results, rather than the expected page for ratings and user
1205 reviews. In other cases, when actions like filling out forms are required, the greedy nature of LLMs
1206 leads them to input more task-relevant information than necessary. This results in a narrower range
1207 of information that can be extracted from the webpage, as shown in the second line in Table 14.
1208 Meanwhile, the limitations of browser automation tools currently prevent the complete restoration of
1209 a web page to its state before action execution. Memory management of web agents also could not
1210 eliminate the effect of past incorrect trajectories. These all highlight the challenges of autonomous
1211 agent reasoning.

1212

1212 G.2 PREMATURE TERMINATION OF TASKS

1213

1214 In the experiments, we also discovered that the web agent sometimes only partially completes tasks.
1215 This typically indicates that web agent sometimes prematurely judges itself as having finished the
1216 task. The reasons for premature termination are varied. For instance, the agent might hallucinate
1217 during inference (such as simplifying a task of reaching a page and filling out content to just reaching
1218 the page), leading it to self-judge the task as complete after only finishing intermediate steps and not
1219 continuing further. In other instances, it may have the right thought process in earlier steps, but fails
1220 to deliver the correct action input or effectively execute the action on the page, yet in subsequent
1221 steps, it “reads” this thought and mistakenly believes the action has been executed. Lastly, when it is
1222 difficult to continue along the current path, the agent might lower its standards for task completion
1223 and erroneously judge the task as complete, thus terminating the task prematurely. As shown in the
1224 third line in Table 14, in the task “Track the status of a rebate for Monty Lue, house number 4847, zip
1225 code 10019 in Menards”, the web agent reached the “Track Your Rebate” page but did not continue
1226 to complete the form, instead prematurely deciding the task was complete and ending the task.

1227

1227 G.3 INFORMATION LOSS IN OBSERVATION

1228

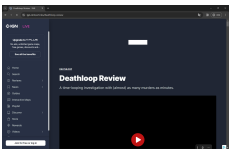
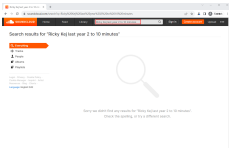
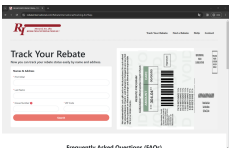
1229 The relationships between web elements are varied and complex. Often, the essential information of
1230 an element is not contained within the element itself but is instead found within its child elements,
1231 parent, or even sibling elements. For instance, a button tag might not always contain useful attributes;
1232 sometimes, they are empty or irrelevant. Based on our understanding of the DOM tree on the web,
1233 we map information from specific elements (like span) to interactive elements such as buttons. Due
1234 to the diversity of these mapping relationships, our framework currently only considers mapping
1235 valuable information from certain special elements to their parent elements, recursively iterating until
1236 an interactive element is identified, as shown in Figure 15. If this recursive search fails to find an
1237 interactive element or reaches the recursion limit, the element is discarded, as illustrated in Figure
1238 16. Given the complexity of webpage elements, our initial implementations focus predominantly on
1239 parent-child mapping relationships. Future work will delve deeper into inter-element mappings to
1240 ensure the accuracy and correctness of element mappings.

1241

1241

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Table 14: Case study of failure trajectories.

State	Task Instruction	Agent's Thought
	Check the rating and user reviews for the game "Deathloop" in ign	To view the full review and user comments for 'Deathloop' on IGN, I need to click on the IGN review link provided in the search results.
	Find a Ricky Kej track to listen and share which has been added in the last year and is between 2 to 10 minutes in soundcloud	To find a Ricky Kej track that was added in the last year and is between 2 to 10 minutes long on SoundCloud, I need to utilize the search and filter features. Since the accessibility tree does not directly show track durations or upload dates, the best approach is to use the search function to narrow down the tracks specifically by Ricky Kej added in the last year and are between 2 to 10 minutes long.
	Track the status of a rebate for Monty Lue, house number 4847, zip code 10019 in menards	After navigating through various steps and overcoming challenges related to the website's interaction model, we've successfully reached the point where the 'Track Your Rebate' form is ready for Monty Lue's information. The task is marked as 'finished' because the user has been guided to the correct section for submission, aligning with the goal of tracking the rebate status.

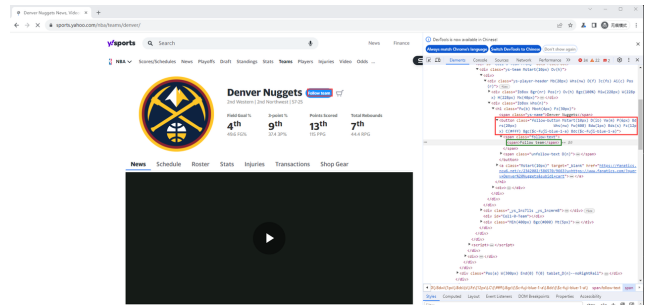


Figure 15: Example on parent-child element mapping strategy

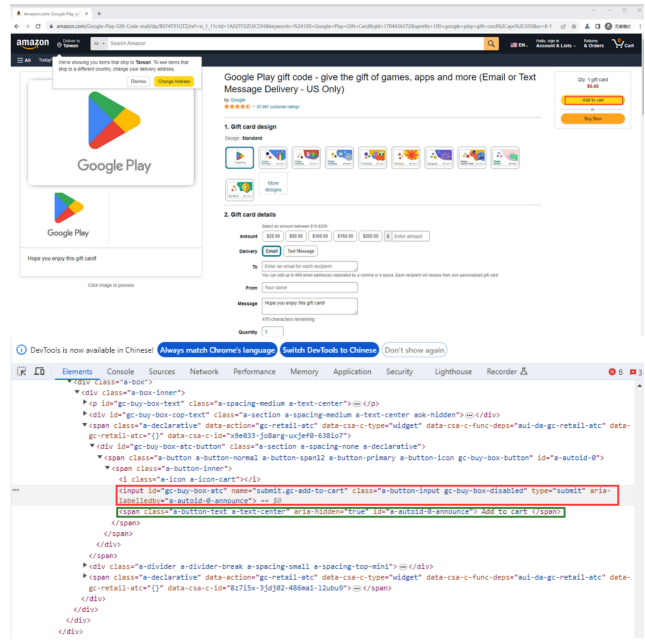


Figure 16: Example on failure case of parent-child element mapping strategy

H DATA VALIDITY TEST REPORT

See Figure 17.

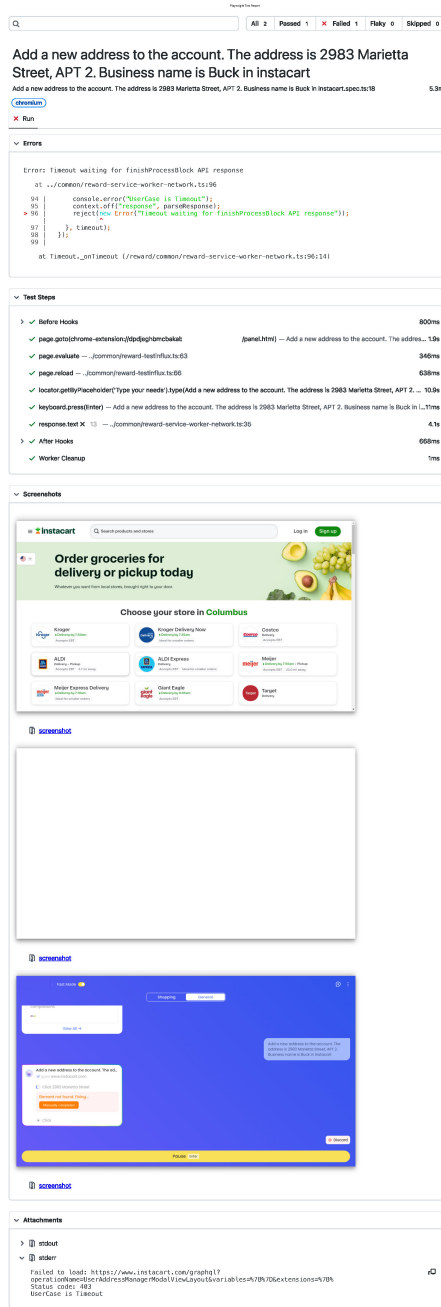


Figure 17: Data validity test report

I EXAMPLES OF MORE ANNOTATED SAMPLES

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

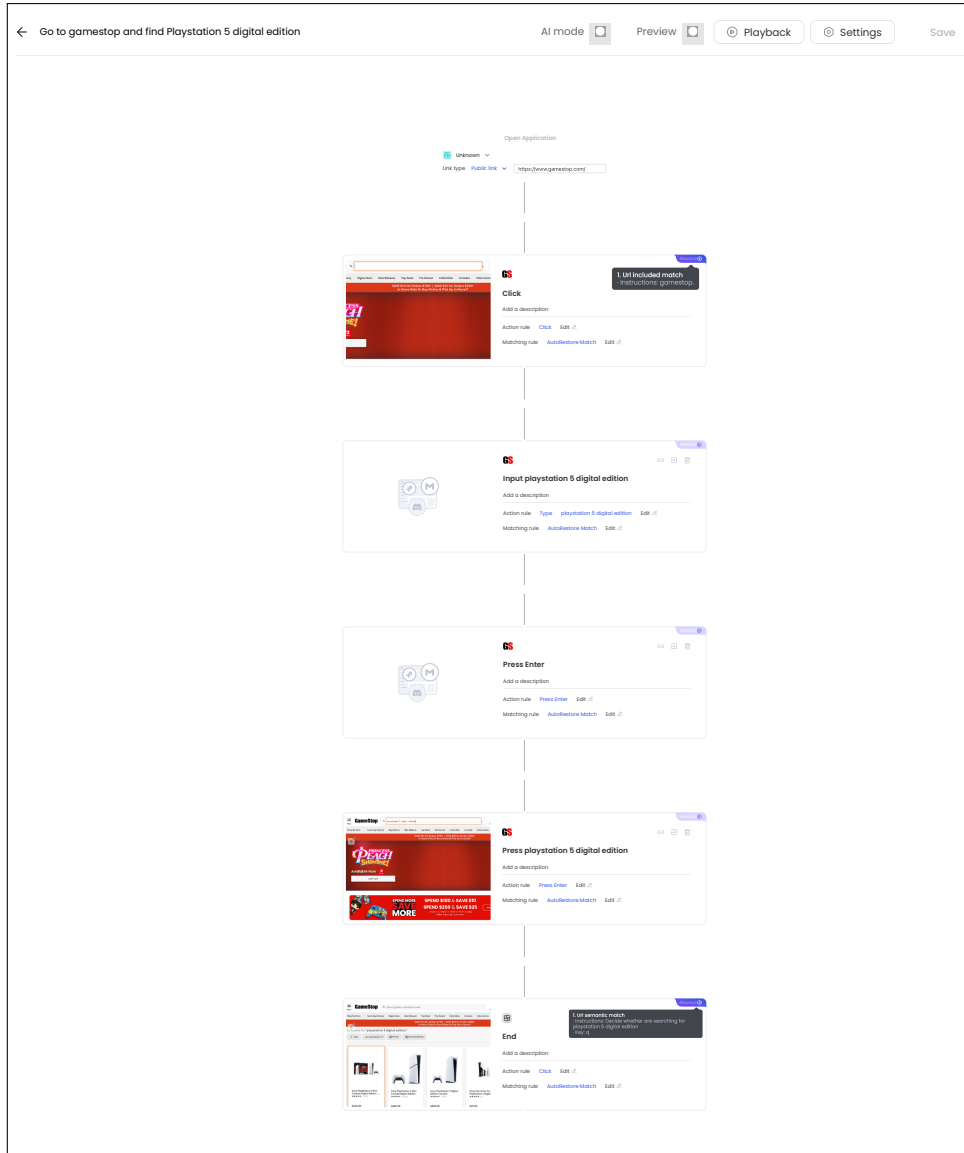


Figure 18: Example on the annotated interface and evaluation function for the task “Go to gamestop and find PlayStation 5 digital edition”

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

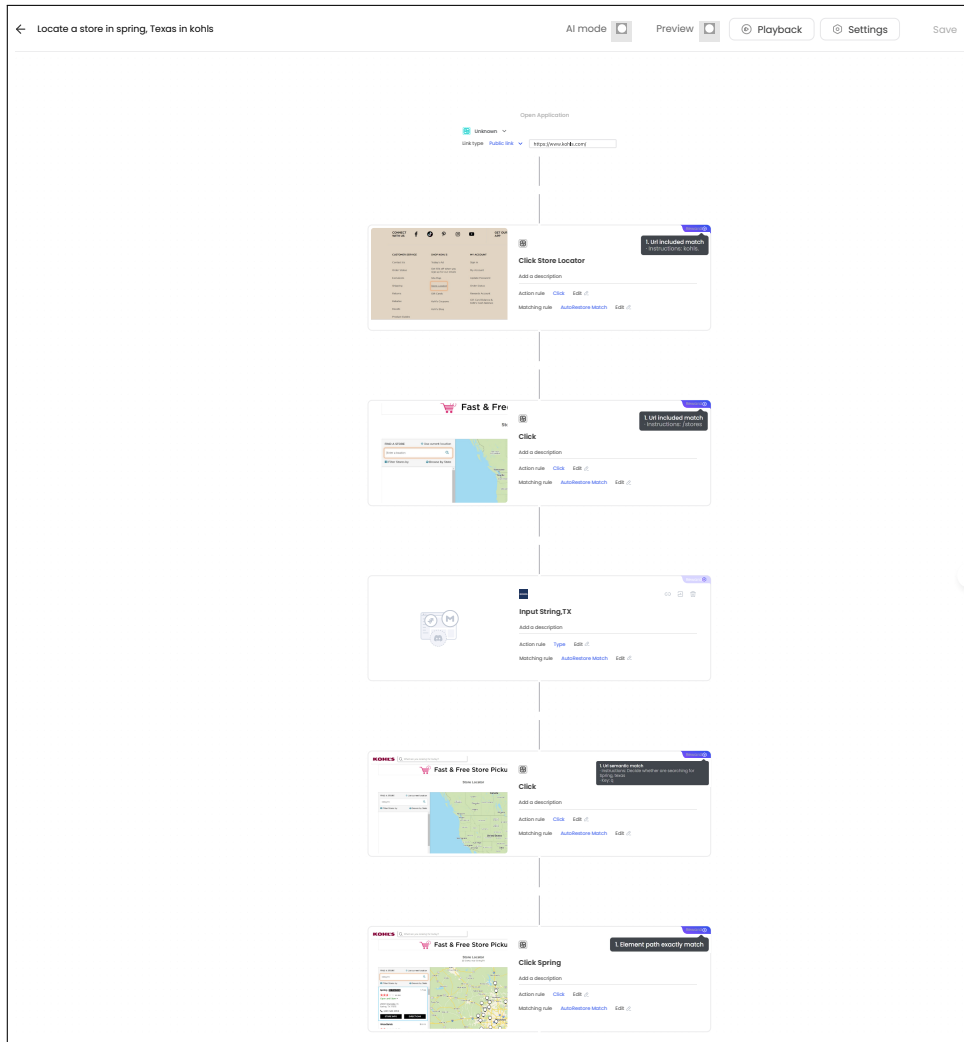


Figure 19: Example on the annotated interface and evaluation function for the task “Locate a store in spring, Texas in kohls”

1458 J LIMITATIONS & FUTURE WORKS

1459

1460

1461

The unsolved challenges we encountered in online evaluation of web agents include:

1462

1463

1464

1465

1466

1467

1468

1. Network Instability: The variability in network conditions can lead to discrepancies between the results obtained from online real-time evaluations and those from closed environments. For instance, issues such as CAPTCHAs, network outages, or inconsistencies across different IPs can influence outcomes. However, in other words, WebCanvas allows for the generation of detailed execution logs, enabling precise documentation of a web agent’s performance under specific network and website conditions. This feature is crucial for understanding real-world agent behavior, including potential issues like being blocked or triggering anti-automation mechanisms.

1469

1470

1471

1472

1473

1474

2. Complex Task Pathways: The diversity of potential execution paths for a given task may not be completely identified by human annotators. This oversight can lead to a misalignment between the defined key nodes and the essential components of task completion, inadvertently penalizing correct processes. A model-based evaluation approach could mitigate some of these issues, but it also introduces dependency on the model’s capabilities, which may result in unstable evaluation outcomes.

1475

1476

1477

1478

1479

1480

3. Static Evaluation Functions: The current static nature of our evaluation functions does not accommodate changes in task instructions based on environmental variables such as time, location, or weather conditions. For example, a task might involve booking a flight to Hawaii next month if the weather is favorable. Ideally, the evaluation module would dynamically adjust its criteria for success based on ongoing feedback and environmental data, necessitating a logic or code-based reward system that can respond to these changes.

1481

1482

1483

1484

1485

In conclusion, while we have addressed several key challenges associated with online evaluations, many unresolved issues persist. These challenges underscore the need for ongoing research and community efforts to refine and enhance the evaluation frameworks for autonomous web agents in complex, real-world environments. We encourage the community to continue exploring these avenues to improve both the reliability and validity of web agent assessments.

1486

1487

1488

K PROMPTS OF PLANNING AND REWARD MODULE

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

Planning Prompt

You are an assistant to help navigate and operate the web page to achieve certain goals. Answer the following questions as best as you can.

There are key information you will get:

****Key Information**:**

- Previous trace: all thoughts, actions and reflections you have made historically.
- Accessibility tree: characteristic expression of the current web page.

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

****Introduction to Accessibility Tree**:**

The accessibility tree is a tree-like data structure that describes the relationships between elements on a web page and provides accessibility information for each element (such as text, links, form elements, etc.).

- ****Accessibility Tree Example**:**

Here is an example of an accessibility tree:

```
```
```

```
current web tab name is 'Google'
```

```
[40] link 'About'
```

```
[41] link 'Store'
```

```
[186] link 'Gmail'
```

```
[187] link 'Images'
```

```
[163] textarea 'Search'
```

```

1512 [236] button 'See more'
1513 '''
1514
1515 In this example, each row represents the characteristic
1516 representation of a web page element. It has three attributes:
1517 '[40]' for the element's element_id, 'link' indicates the element
1518 is a link, and 'About' for the content of the element.
1519 Note: The above element provided is purely for illustrative
1520 purposes and should NEVER be used directly in your output!
1521
1522 You should always consider previous and subsequent steps and what
1523 to do.
1524 **Thought Space**:
1525 - What action do you think is needed now to complete the task?
1526 - What's the reason of taking that action?
1527
1528 You have access to the following tools(helpful to interact with web
1529 page):
1530 **Execution Action Space**:
1531 - goto: useful for when you need visit a new link or a website,
1532 it will open a new tab.
1533 - fill_form: useful for when you need to fill out a form or
1534 input something from accessibility tree. Input should be a
1535 string.
1536 - google_search: useful for when you need to use google to
1537 search something.
1538 - click: useful for when you need to click a button/link from
1539 accessibility tree.
1540 - select_option: useful for when you need to select a drop-down
1541 box value. When you get (select and option) tags from the
1542 accessibility tree, you need to select the serial number(
1543 element_id) corresponding to the select tag, not the option,
1544 and select the most likely content corresponding to the option
1545 as Input.
1546 - go_back: useful when you find the current web page encounter
1547 some network error or you think the last step is not helpful.
1548
1549 You also need to provide an effective description of the current
1550 execution action.
1551 A proper description contains:
1552 - What website it is;
1553 - Which action you choose;
1554 - REMEMBER DO NOT LEAVE THE DESCRIPTION EMPTY!
1555
1556 You have to follow the instructions or notes:
1557 **Important Notes**:
1558 - Under the following conditions, you are restricted to using
1559 the 'google_search' or 'goto' tools exclusively:
1560 1. In the initial step of a process or when there's no
1561 preceding interaction history (i.e., the previous trace is
1562 empty).
1563 2. In situations where the accessibility tree is absent or
1564 not provided.
1565 - Your action should not be the same as last step's action.
1566 - The 'element_id' should be an integer accurately representing
1567 the element's ID in the accessibility tree.
1568 - AVOID using the provided example's element_id as your output.
1569 - The output JSON-formatted code block must be valid; otherwise
1570 , it cannot be recognized.
1571
1572 **Special Circumstances Guidelines**:
1573 - When performing a search on a website, if you find the search
1574 results do not display sufficient content, consider
1575 simplifying or modifying your search query. Reducing the

```

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

complexity of your search query or altering keywords may yield more comprehensive results.

Please ensure the accuracy of your output, as we will execute subsequent steps based on the 'action', 'action\_input' and 'element\_id' you provide.

**\*\*Output Requirements\*\*:**

- Ensure your output strictly adheres to the JSON-formatted code block outlined below:

```
'''
{
 "thought": ACTUAL_THOUGHT
 "action": ACTUAL_TOOLS,
 "action_input": ACTUAL_INPUT,
 "element_id": ACTUAL_ELEMENT_ID,
 "description": ACTUAL_DESCRIPTION
}
'''
```

- A VALID JSON-FORMATTED CODE BLOCK EXAMPLE AS FOLLOWS:

```
'''
{
 "thought": "In order to complete this task, I need to go to
 the Google home page",
 "action": "click",
 "action_input": "button",
 "element_id": "236",
 "description": "Now I\'m on Google\'s main page. I\'m now
 clicking the button with element_id [236] to see more
 information."
}
'''
```

### Reward Prompt

You are an assistant to help navigate and operate the web page to achieve certain task.  
Your goal is to evaluate the previous series of traces(thoughts and actions) and think about what key steps are needed to complete the task in the future.

There are key information you will get:

**\*\*Key Information\*\*:**

- Previous trace: all thoughts, actions and reflections you have made historically.
- Accessibility tree: characteristic expression of the current web page.
- Screenshot: visual information of the current web page (may include).

You also need to combine the previous trace to give the completion status of the current task.

**\*\*Status Of Task Completion\*\***

- doing: You have completed the intermediate steps of the target task but not entirely finish the target task.
- finished: You are entirely certain about completing the target task.
- loop: You find that the the last two steps of previous actions are the same, it is determined that the process is stuck in a local optimum solution.

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

You will judge and score the task completion and reasonableness of previous actions. The score ranges from 1-10, but the score you give can only be selected from [1, 3, 7, 9, 10].

**\*\*Judging and Scoring Criteria\*\*:**

- score = 1: You find that the status of the task is stuck in a loop by analyzing the previous trace.
- score = 3: You find that performing the previous trajectories (thoughts and actions) is not likely helpful in completing target task and you need to adjust the direction of your planning and action or start over from beginning.
- score = 7: You find that performing the previous trajectories (thoughts and actions) are helpful in completing the target task.
- score = 9: You find that performing the previous trajectories (thoughts and actions) are a very critical intermediate step to complete this task.
- score = 10: You find that performing the previous trajectories (thoughts and actions) have completed the task perfectly.

You need to provide an effective evidence of scoring for the series of the previous trace.

- Why do you give this score?
- What is the reason?

You also need to provide an effective description or summary of the above requirements through key information and characteristics of the current web page.

**\*\*A proper description contains\*\*:**

- What is the current completion status of the task? (IMPORTANT)
- REMEMBER DO NOT LEAVE THE DESCRIPTION EMPTY!

**\*\*Output Requirements\*\*:**

- Ensure your output strictly follows this format:

```
```json
{
  "status": "ACTUAL_STATUS",
  "score": "ACTUAL_SCORE",
  "reason": "ACTUAL_REASON",
  "description": "ACTUAL_DESCRIPTION"
}
```

- A VALID JSON-FORMATTED CODE BLOCK EXAMPLE AS FOLLOWS:

```
```
{
 "status": "doing",
 "score": "3",
 "reason": "You need to complete a search for camping tents that can accommodate 2 people and sort the results in rei by price from low to high. According to your previous trajectory, you navigated to the rei official website and clicked the 2-person button, which are correct actions. But when you complete the final step of sorting prices, you actually click on a link to a tent product. This is a completely unreasonable action. So I give it 3 points."
 "description": "According to the current web page information, you can know that this is the homepage of a tent product, which is not very consistent with the purpose of the target task."
}
```

1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

### Reward Prompt - With Golden Reference

You are an assistant to help navigate and operate the web page to achieve certain task.  
Your goal is to evaluate the previous series of traces(thoughts and actions) and think about what key steps are needed to complete the task in the future.

There are key information you will get:

**\*\*Key Information\*\*:**

- Previous trace: all thoughts, actions and reflections you have made historically.
- Current Webpage Information:
  - Accessibility tree: characteristic expression of the current web page.
  - Screenshot: visual information of the current web page. ( may include)
- Reference Guide: detailed and step-by-step reference guide for completing the target task, serving as a benchmark for evaluating progress and strategizing the necessary actions.

**\*\*Notes to Reference Guide\*\*:**

- The Reference Guide plays a crucial role in aiding the evaluation of the current Status of Task Completion. The ' Completion Verification' section within the Reference Guide is instrumental in determining whether a task can be classified as 'finished.'
- Furthermore, for a task to be considered fully completed, all **\*\*key conditions\*\*** must be met as specified.

You also need to combine the previous trace to give the completion status of the current task.

**\*\*Status of Task Completion\*\***

- doing: You have completed the intermediate steps of the target task but not entirely finish the target task.
- finished: You are entirely certain about completing the target task.
- loop: You find that the the last two steps of previous actions are the same, it is determined that the process is stuck in a local optimum solution.

You will judge and score the task completion and reasonableness of previous actions. The score ranges from 1-10, but the score you give can only be selected from [1, 3, 7, 9, 10].

**\*\*Judging and Scoring Criteria\*\*:**

- score = 1: You find that the status of the task is stuck in a loop by analyzing the previous trace.
- score = 3: You find that performing the previous trajectories (thoughts and actions) is not likely helpful in completing target task and you need to adjust the direction of your planning and action or start over from beginning.
- score = 7: You find that performing the previous trajectories (thoughts and actions) are helpful in completing the target task.
- score = 9: You find that performing the previous trajectories (thoughts and actions) are a very critical intermediate step to complete this task.
- score = 10: You find that performing the previous trajectories(thoughts and actions) have completed the task perfectly.

You need to provide an effective evidence of scoring for the series of the previous trace.

- Why do you give this score?



```

1728 - What is the reason?
1729
1730 You also need to provide an effective description or summary of the
1731 above requirements through key information and characteristics of
1732 the current web page.
1733 **A proper description contains**:
1734 - What is the current completion status of the task? (IMPORTNAT
1735)
1736 - REMEMBER DO NOT LEAVE THE DESCRIPTION EMPTY!
1737
1738 **Output Requirements**:
1739 - Ensure your output strictly follows this format:
1740 ```json
1741 {
1742 "status": "ACTUAL_STATUS",
1743 "score": "ACTUAL_SCORE",
1744 "reason": "ACTUAL_REASON",
1745 "description": "ACTUAL_DESCRIPTION"
1746 }
1747 ```
1748 - A VALID JSON-FORMATTED CODE BLOCK EXAMPLE AS FOLLOWS:
1749 ```
1750 {
1751 "status": "doing",
1752 "score": "3",
1753 "reason": "You need to complete a search for camping tents
1754 that can accommodate 2 people and sort the results in rei
1755 by price from low to high. According to your previous
1756 trajectory, you navigated to the rei official website and
1757 clicked the 2-person button, which are correct actions.
1758 But when you complete the final step of sorting prices,
1759 you actually click on a link to a tent product. This is a
1760 completely unreasonable action. So I give it 3 points."
1761 "description": "According to the current web page
1762 information, you can know that this is the homepage of a
1763 tent product, which is not very consistent with the
1764 purpose of the target task."
1765 }
1766 ```

```

### Semantic Match Prompt

```

1766 Now you are an assistant to judge whether 2 elements are
1767 semantically same. I'll provide a judge rule and an answer.
1768 If they are the same, you should return 1. If they are not related,
1769 you should return 0.
1770 If they are related but not identical, return a decimal (two
1771 decimal places) between 0 and 1 of the degree of relevance you
1772 think.
1773 For example, the judge rule is: Decide whether the place is New
1774 York. The score of "new york" and "New York" are both 1, "Brooklyn
1775 " should be 0.
1776 However, if the judge rule is: Decide whether the place is in New
1777 York. The score of "new york" and "New York" and "Brooklyn" are
1778 all 1.
1779 Another example, the judge rule is: Decide whether I'm looking for
1780 clothes. The score of "red Clothes" and "green jacket" should also
1781 be 1.
1782 However, if the judge rule is: Decide whether I'm looking for red
1783 clothes. the score of "bright red Clothing" could be 0.85(red

```

1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835

include bright red but they are not the same), the score of "green  
Clothes" should be 0.5 (red is not green).  
Remember, you should return a number with " and an explanation.  
Like output: "1", (your explanation)