
WebGraphEval: Multi-Turn Trajectory Evaluation for Web Agents using Graph Representation

Yaoyao Qian*
Northeastern University

Yuanli Wang
Boston University

Jinda Zhang
University of Victoria

Yun Zong
University of Minnesota

Meixu Chen
Northeastern University

Hanhan Zhou
George Washington University



Jindan Huang
Tufts University

Yifan Zeng
Oregon State University

Xinyu Hu
University of Texas at San Antonio

Chan Hee Song
The Ohio State University

Danqing Zhang
PathOnAI.org

 Project Website  Demo Website

Abstract

Current evaluation of web agents largely reduces to binary success metrics or conformity to a single reference trajectory, ignoring the structural diversity present in benchmark datasets. We present **WebGraphEval**, a framework that abstracts trajectories from multiple agents into a unified, weighted action graph. This graph representation is directly compatible with existing benchmarks such as WebArena, using both leaderboard trajectories and newly collected runs, and provides a principled basis for analyzing solution spaces without modifying environments. The framework canonically encodes actions, merges recurring behaviors, and applies structural analyses including reward propagation and success-weighted edge statistics. Evaluations across thousands of trajectories from six web agents demonstrate that the graph abstraction captures cross-model regularities, highlights redundancy and inefficiency, and identifies critical decision points overlooked by outcome-based metrics. By framing web interaction as graph-structured data, WebGraphEval establishes a general methodology for multi-path, cross-agent, and efficiency-aware evaluation of web agents.

1 Introduction

Benchmarks such as WebArena [1] and Mind2Web [2] provide thousands of recorded trajectories showing how web agents interact with interfaces. These datasets contain many valid solution paths for the same task, but current evaluation methods reduce performance to binary success or to matching a reference trajectory. This leaves out the information in the trajectories themselves, including how agents explore, take detours, or recover.

Existing approaches show this limitation in different forms. Outcome-based metrics ignore the intermediate process. The "LLM-as-a-Judge" paradigm [3] allows more flexible judgments but still reduces evaluation to a final outcome. Trajectory conformity [4, 2] compares against a reference path but reflects the bias of the chosen path, penalizes alternative strategies, and requires frequent updates

*Corresponding author. qian.ya@northeastern.edu

Table 1: Comparison of Web Agent Evaluation Methods

Method/Benchmark	Graph-based Analysis	Multi-path Support	Analysis Scope	Judge Method
WebShop [5]	✗	✗	Final message (Single)	Rule-based
Mind2Web [2]	✗	✓	Single Trajectory	Rule-based
WebArena [1]	✗	✗	Single Trajectory	Rule-based + LLM-as-judge
VisualWebArena [6]	✗	✗	Final message (Single)	Rule-based + LLM-as-judge
WebVoyager [7]	✗	✗	Final message (Single)	LLM-as-judge
Mind2Web 2.0 [8]	✗	✗	Single Trajectory	Agent-as-judge
AgentBoard [9]	✓	✗	Single Trajectory	LLM-as-judge
ST-WebAgentBench [10]	✗	✗	Single Trajectory	Rule-based
VideoWebArena [11]	✗	✗	Single Trajectory	Rule-based + LLM-as-judge
WebGraphEval (Ours)	✓	✓	Cross Multi Trajectories	LLM-as-judge

as interfaces change. In all cases, the common issue is that multiple trajectories are collected, but their overlap, differences, and repeated errors are not used.

We propose that evaluation should use a structured representation of trajectories rather than isolated sequences. Graphs are a natural choice: nodes represent actions, edges represent transitions, and multiple trajectories can be combined into a single graph. This makes it possible to see common strategies, points where agents fail, and how different models behave on the same task.

We present **WebGraphEval**, a framework that builds such graphs from raw trajectories. The framework converts actions into a canonical form, merges similar behaviors into weighted nodes and edges, and applies analyses such as reward propagation and success-weighted edge classification. Using leaderboard data and additional runs from WebArena, WebGraphEval works with existing benchmarks and can be extended to new agents. Across six web agents, it shows how graph-based evaluation captures shared strategies, redundant actions, and critical steps that outcome-based metrics cannot measure. This work makes four contributions to web agent evaluation:

1. **WebGraphEval framework.** We present WebGraphEval, a graph-based evaluation framework that aggregates multiple trajectories into a weighted action graph. This representation captures both shared strategies and divergent behaviors, enabling structured analysis that outcome- or conformity-based methods overlook.
2. **Cross-agent and benchmark compatibility.** WebGraphEval is applied to leaderboard trajectories and additional runs from WebArena, covering six different agents. This shows that the method works directly with existing benchmarks and enables systematic cross-model comparison.
3. **Analytical methods.** The framework includes canonicalization of actions, node merging, and structural analyses such as reward propagation and success-weighted edge statistics. These methods allow us to study task difficulty, agent efficiency, and decision points in a unified way.
4. **Practical protocol and tool.** We design an LLM-based annotation protocol to assign necessity labels with 78% agreement against human judgments, making graph construction scalable. We also provide a visualization interface (<https://web-graph-eval.vercel.app/>) for exploring trajectories before and after graph construction.

2 Related Work

Benchmarks and Evaluation Protocols Benchmarks for web agents range from simulated settings (WebShop [5]) to live websites (Mind2Web [2], WebArena [1]), with extensions to vision, cross-domain, and safety tasks [6, 11, 12, 10]. Most leaderboards report binary success. Other metrics include progress rate (AgentBoard [9]), skill scoring (FLASK [13]), step-wise conformity to a reference [4, 2], and LLM-as-a-Judge outcome checks [8]. These protocols evaluate each trajectory in isolation and do not analyze the structure across multiple valid paths.

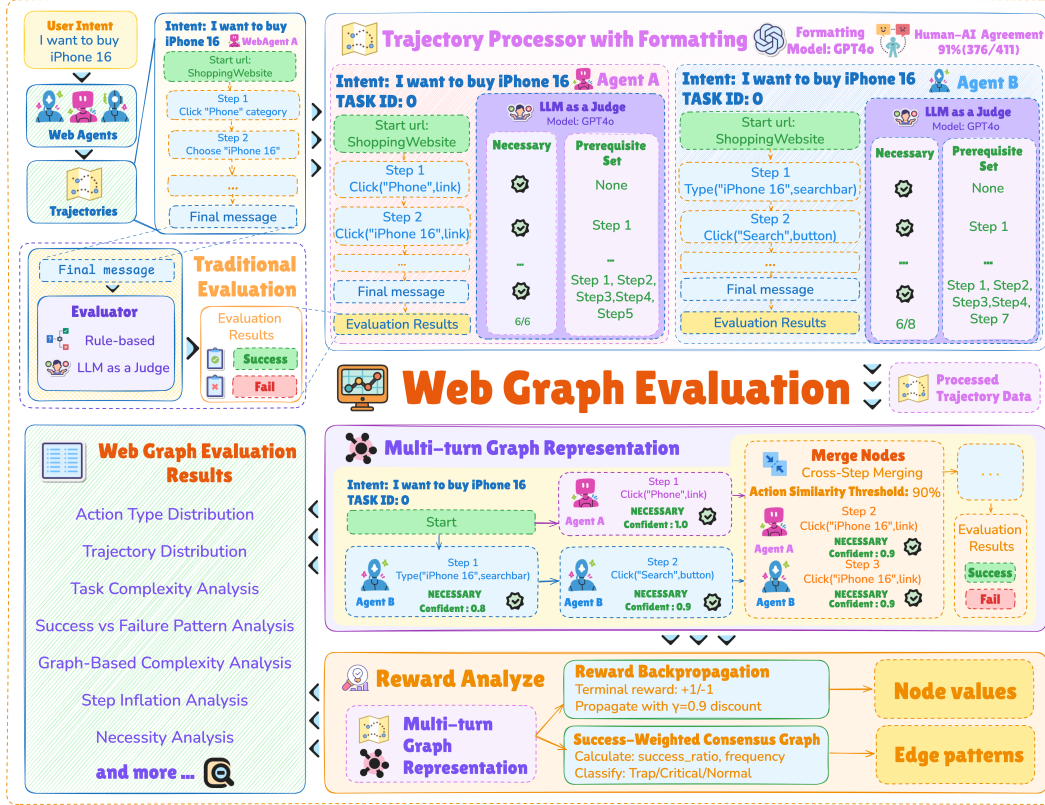


Figure 1: Overview of the WebGraphEval pipeline. Trajectories from multiple agents are first pre-processed and canonicalized into standardized action units. These are merged into a directed action graph, where nodes represent actions and edges represent observed transitions. The graph is then analyzed with reward backpropagation and success-weighted edge classification, and the resulting structure supports multi-dimensional evaluation of agent efficiency, redundancy, and strategy overlap.

Graph-Based and Structural Approaches Graphs have been applied in many domains to represent sequential decisions and relations. In embodied AI, graphs are used for task planning [14]. In retrieval-augmented generation, systems such as GraphRAG [15], LightRAG [16], and Causal GraphRAG [17] build graphs over entities or events to support multi-hop reasoning in document or news QA. Graph-R1 [18] uses a hypergraph structure for complex QA tasks. For web agents, WebVoyager [7] analyzes navigation patterns and Go-Browse [19] frames exploration as graph search. These works show that graphs can support reasoning and analysis, but they have not been used as a systematic evaluation framework. A side-by-side summary of these approaches is shown in Table 5. In contrast, WebGraphEval builds directed graphs over canonicalized actions and transitions, aggregates trajectories across agents, and applies graph analysis for evaluation.

3 Methodology

WebGraphEval is a graph-based evaluation framework that aggregates trajectories from multiple agents into a unified structure (Figure 1). Web interactions are represented as directed transitions between actions: nodes correspond to canonicalized actions, and edges represent observed transitions weighted by frequency and outcome. This abstraction preserves both common strategies and divergent behaviors, enabling evaluation beyond single-path outcomes.

The framework integrates two complementary analysis methods. Reward backpropagation propagates task outcomes backward through trajectories to estimate the value of earlier actions, while success-weighted edge classification labels transitions by their empirical association with success or failure. These provide temporal and structural perspectives on agent behavior.

As illustrated in Figure 1, the pipeline unfolds in four stages. Trajectories are first pre-processed and canonicalized into standardized actions. Canonical actions are then merged into nodes, and transitions are stored as directed edges. The resulting graph is analyzed with reward backpropagation and edge classification. Finally, the graph supports multi-dimensional evaluation of agents, including efficiency, redundancy, and strategy comparison across models.

3.1 Trajectory Formalism and Pre-processing

The first stage of WebGraphEval is to convert raw agent trajectories into a formal representation suitable for analysis. Each trajectory $\tau = (a_1, \dots, a_n)$ is an ordered sequence of actions, where each action $a_i = \langle \text{desc}_i, \text{url}_i \rangle$ consists of a textual description and the corresponding execution URL. A key challenge is that agents often describe the same action in different natural language forms. To address this, we use an LLM-based canonicalization step that maps diverse descriptions to a standardized schema. For example, “clicked on the ‘Submit’ button” and “press Submit” are both converted to the unified form `click(text='Submit', element='button')`.

Once actions are normalized, trajectories are labeled as successful or failed using an LLM-as-Judge protocol. The model is prompted with task-specific rules and evaluates whether the final state of each trajectory satisfies the original task requirements. This step produces two disjoint sets: successful trajectories D_{success} and failed trajectories D_{fail} .

Finally, we refine the representation by identifying redundant actions. Each action a_i is assigned a necessity label $\nu(a_i) \in \{0, 1\}$, where $\nu(a_i) = 1$ indicates that the action is essential for task completion. This annotation is generated automatically by an LLM and spot-checked with human evaluators to confirm reliability. Figures 13 and 14 show the prompts used for action canonicalization and necessity labeling.

3.2 Consensus Graph Construction

Given canonicalized trajectories, we construct a consensus directed graph $G = (V, E)$ that serves as the representation for subsequent analysis. The input is a set of action sequences $\tau = (a_1, \dots, a_n)$ from multiple agents, and the output is a graph in which each node $v \in V$ corresponds to an equivalence class of actions while each edge $(u \rightarrow v) \in E$ represents an observed transition. Nodes and edges carry counts and outcome-conditioned statistics that are used in later stages of evaluation.

To create the node set, we merge semantically similar actions into common nodes. Each action is already expressed in a structured format after canonicalization, such as `click(text=..., element=...)`. We measure the similarity between two actions a_i and a_j with a normalized edit distance,

$$\text{sim}(a_i, a_j) = 1 - \frac{\text{Levenshtein}(a_i, a_j)}{\max(|a_i|, |a_j|)}.$$

Merging is applied in two passes. In step-based merging, actions are only compared with others that occur at the same time index across trajectories, which preserves temporal roles. In cross-step merging, comparisons extend across different positions, which allows the graph to capture actions that recur in multiple places. The merging process is implemented with a deterministic union-find procedure to avoid dependence on processing order.

The similarity threshold θ is a tunable parameter. A high value leaves near-duplicates unmerged, producing fragmented graphs, while a low value risks collapsing distinct actions. We use $\theta = 0.9$ as a practical compromise: it removes superficial wording differences while still maintaining distinctions between genuinely different behaviors. This threshold works without task-specific tuning, although future work could replace it with embedding-based similarity or learned clustering.

Edges are then created from adjacent actions in each trajectory. For every pair (a_t, a_{t+1}) , we add a directed edge $(\pi(a_t) \rightarrow \pi(a_{t+1}))$, where π maps actions to their merged node. Each edge stores the number of times it occurs, its counts under successful and failed trajectories, and the empirical success rate derived from these counts. Nodes also retain step-index histograms to support later analysis of temporal patterns. The pairwise similarity computation has quadratic complexity in the number of actions. Since our dataset is moderate in size, we perform the comparisons directly. For larger datasets, approximate blocking strategies (e.g., by action type or host) could reduce the cost to near-linear, but we leave this as future work.

3.3 Dual Reward Mechanisms

The consensus graph aggregates trajectories from all six agents, but to make sense of this structure we need a way to measure which actions and transitions matter most. WebGraphEval employs two complementary reward mechanisms. The first propagates outcome signals along entire trajectories, highlighting how early decisions contribute to eventual success or failure. The second examines transition statistics directly, identifying structural patterns such as traps, bottlenecks, and critical paths. Together these methods turn raw trajectory data into interpretable evidence about agent behavior.

The temporal reward backpropagation mechanism is inspired by reinforcement learning. Each terminal state is initialized with a reward of +1 if the task succeeds and -1 if it fails. These values are then propagated backward through the graph so that earlier nodes inherit credit or blame from their successors. Formally, the value $V(v)$ of a node v is defined in terms of its outgoing edges: each successor contributes proportionally to how often the transition is observed, discounted by a factor $\gamma = 0.9$ to reduce the influence of distant outcomes. This procedure assigns positive values to actions that reliably lead, through several steps, to task completion, and negative values to those that tend to end in failure. In this way, the graph reveals decision points that matter even if their effects only appear later in the trajectory.

The success-weighted consensus analysis takes a different perspective. Instead of propagating signals through paths, it looks at the observed frequency and success ratio of each edge. For an edge e , the success ratio $s(e)$ is the proportion of times it appears in successful trajectories, and the weight $w(e)$ is its relative frequency compared to all edges. Using these two measures, edges are classified into interpretable categories. Trap edges are those that occur often but almost always lead to failure, such as clicking a misleading link. Critical edges are rare but consistently successful transitions, representing expert-like behavior. Bottleneck edges are high-frequency transitions with mixed success, often corresponding to fragile but necessary steps like form submissions. All other edges are treated as normal, reflecting routine navigation without strong outcome bias. This classification exposes how agents collectively approach tasks and where they are most prone to error.

Beyond individual edges, we also estimate the importance of each node. A node is considered important if it is both frequently visited and strongly associated with successful outcomes. We capture this with a score that averages the success ratios of incoming and outgoing edges and multiplies by the node’s visitation frequency. Intuitively, this highlights actions that not only appear often in trajectories but also play a decisive role in task completion. An example would be reaching a login page, which is both common and essential for downstream success.

3.4 Multi-dimensional Behavioral and Cross-Agent Evaluation

The final stage of WebGraphEval evaluates agent behavior across both individual and comparative dimensions. At the trajectory level, we assess path optimality by comparing observed trajectory lengths with the shortest successful paths in the dataset, excluding anomalous cases. This metric captures inefficiency arising from unnecessary exploration. Leveraging LLM-based annotations, we further distinguish necessary from redundant actions, providing a fine-grained measure of how directly agents pursue task goals. Temporal dynamics are also incorporated by analyzing how action types are distributed across early, middle, and late phases of a trajectory, revealing when redundancy or recovery attempts are most likely to occur.

These trajectory-level measurements are then aggregated into agent-level profiles. Beyond raw success rates, we examine performance across task categories, complexity levels, and trajectory lengths, which highlight characteristic strategy patterns such as average path length, preferred navigation choices, and resilience to increasing task difficulty. To extend the analysis further, WebGraphEval enables cross-agent comparisons on shared tasks. Using entropy-based metrics and clustering over trajectory distributions, we quantify the diversity of strategies, determining whether models converge on similar behaviors or adopt distinct approaches. The consensus graph also identifies action sequences that consistently appear in successful trajectories across multiple frameworks, pointing to robust strategies that generalize beyond individual agents.

Table 2: Framework-level performance comparison evaluated using llm_judge (single evaluation).

Framework	Success	Failure	Success Rate	Avg Steps	Avg. Confidence	Necessity Rate
Zeta Labs Jace.AI [20]	526	286	64.78%	5.6	0.934	73.7%
IBM CUGA [21]	477	331	59.03%	5.3	0.972	80.6%
Learn by Interact [22]	440	372	54.19%	6.0	0.967	72.9%
UI-TARS [23]	296	468	38.74%	13.2	0.976	82.0%
OpenAI-CUA [24]	226	582	27.97%	6.4	0.932	74.1%
BrowserUse [25]	215	549	28.14%	15.6	0.926	74.7%
Total	2,180	2,588	45.75%	–	–	–

4 Experiments

4.1 Dataset and Setup

Our experiments are conducted on a new trajectory dataset built from the WebArena benchmark [1]. The dataset consists of 4,768 trajectories collected from six agent frameworks attempting 812 unique tasks. The resulting corpus produces large and diverse action graphs, with 40,431 nodes and 45,656 edges in total, averaging 49.79 nodes and 56.23 edges per task. Across all trajectories, we observe 40,888 individual actions, including 19,380 clicks (47.4%), 8,312 type actions (20.3%), and 1,302 select operations (3.2%). Overall, 2,180 trajectories were successful (45.7%) and 2,588 failed (54.3%), providing balanced coverage of both effective and ineffective strategies.

All trajectories were evaluated consistently using an LLM-based judge (o4-mini-2025-04-16, temperature 0.1) to determine success or failure, and action-level necessity labels were obtained with GPT-4o-2024-08-06. Confidence scores assigned by the LLM are consistently high across frameworks (0.926–0.976), indicating stable judgments independent of outcome. To validate annotation quality, each framework was evaluated with three independent runs, and human evaluation was performed on sampled subsets. Canonicalization achieved 91% agreement with human annotators (376/411), while necessity judgments reached 78% agreement (404/520). These checks confirm that LLM-based annotations are sufficiently reliable for downstream analysis.

For graph construction, actions were merged using a normalized edit similarity threshold of $\theta = 0.9$. This threshold reduces superficial variation across agent implementations while preserving distinct behaviors, and ensures that common strategies are captured without erasing meaningful diversity. The resulting graphs reveal that 87.4% of tasks involve at least one merge, though only 5.6% of nodes are merged overall, highlighting both structural overlap and significant diversity in agent behavior.

Finally, two structural characteristics of the dataset are noteworthy. First, a subset of tasks (91 cases, about 13% of all successes) contain one-step successful trajectories, which likely reflect differences in action granularity or pre-filled states; these anomalies are retained for completeness but analyzed separately later. Second, outcome agreement across frameworks is limited: among 761 tasks attempted by all six agents, only 29 tasks (3.8%) were solved universally, 99 tasks (13.0%) failed universally, and the majority (83.2%) showed mixed outcomes. This heterogeneity underscores the need for evaluation methods that integrate across diverse agent behaviors rather than assuming a single reference trajectory.

4.2 Framework-Level Statistics

We first compare performance across the six agent frameworks. Table 2 reports success and failure counts, success rates, average trajectory lengths, confidence scores, and necessity rates. Success rates vary widely, from 64.78% for Jace.AI—the strongest overall performer—to 27.97% for OpenAI-CUA. Trajectory lengths also differ: IBM CUGA completes tasks in an average of 5.3 steps, while BrowserUse requires 15.6. Necessity rates range from 72.9% to 82.0%, reflecting differences in how directly agents pursue task goals. The uniformly high confidence (0.926–0.976) indicates that the LLM-as-judge provides stable and reliable evaluations, ensuring that performance differences reflect true behavioral variation rather than judging noise.

These results demonstrate that efficiency and effectiveness are not tightly coupled. Both UI-TARS and IBM CUGA maintain focused action sequences, as indicated by high necessity rates, yet their success

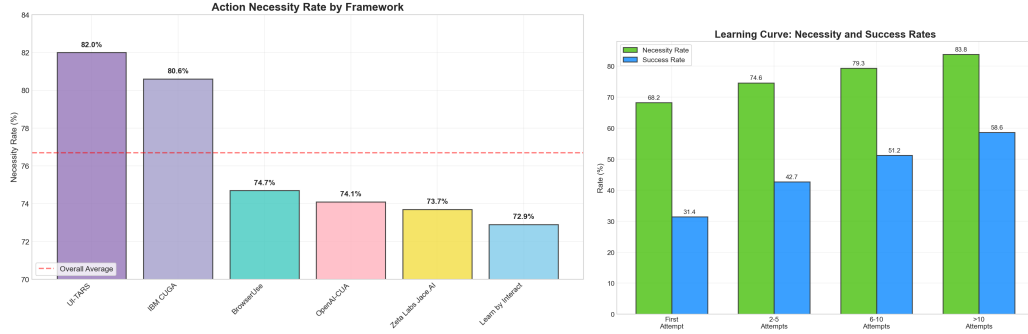


Figure 2: Action necessity across frameworks and learning dynamics. (Left) Necessity rates of six frameworks, showing variation from 72.9% to 82.0%. (Right) Learning curve showing how necessity rates improve with repeated task attempts, suggesting necessity is a learnable signal for agent optimization.

rates differ substantially (38.7% vs. 59.0%). This mismatch suggests that minimizing redundancy alone is insufficient: agents must also make correct decisions at critical points. Figure 2 (left) illustrates this divergence, showing that high necessity does not always translate into high success.

Necessity also evolves with experience. Figure 2 (right) shows a learning curve where necessity rises from 68% on first attempts to over 83% after ten. This trend indicates that necessity is a learnable signal: with repeated exposure, agents reduce redundant actions and become more efficient.

Two structural characteristics further qualify these results. First, 91 tasks (13% of all successful cases) contain one-step successes, likely due to action granularity or pre-filled states. These anomalies are retained but analyzed separately to avoid distorting efficiency metrics. Second, cross-framework agreement is limited: of 761 tasks attempted by all six agents, only 29 (3.8%) were solved universally and 99 (13.0%) failed universally, while the majority (83.2%) showed mixed outcomes. This heterogeneity underscores that no single framework dominates and highlights the value of integrating across diverse agent behaviors rather than relying on a single reference trajectory.

4.3 Graph-Level Structural Analysis

Consensus graphs reveal clear structural differences across agents. IBM CUGA and Learn by Interact, with shorter average trajectories, generate compact and linear graph topologies, whereas UI-TARS and BrowserUse produce highly exploratory graphs with multiple branches. This divergence reflects different balances between goal-directed execution and exploratory search.

Task success follows a non-linear pattern with respect to complexity: simple tasks succeed 44.1%, medium-complexity tasks peak at 54.0%, and performance declines to 45.3% and 31.0% for complex and very complex tasks. We formalize complexity as

$$\text{Complexity} = \frac{\text{nodes} \times \text{edges}}{\text{trajectories}},$$

and observe a 7.9% degradation from low- to high-complexity tasks (Figure 3), confirming structural complexity as a meaningful predictor of difficulty.

Efficiency is further captured by step inflation, the ratio between observed steps and the shortest successful path. The average inflation is 2.14 \times , with some agents taking over 15 steps for 4-step solvable tasks. Frameworks with higher structural efficiency, such as IBM CUGA (80.6% necessity rate), achieve stronger success compared to less focused agents like Learn by Interact (72.9%), showing that necessity and inflation jointly capture key behavioral differences.

As shown in Figure 3 (left), trajectory length and success rates exhibit a clear inverted-U pattern: medium-length trajectories (6–10 steps) achieve the highest success rate (53.4%), while both short (47.6%) and very long trajectories (30.9%) underperform. This indicates that tasks of moderate length provide the right balance between feasibility and planning opportunity, whereas trivial or overly long trajectories expose limitations in current agents.

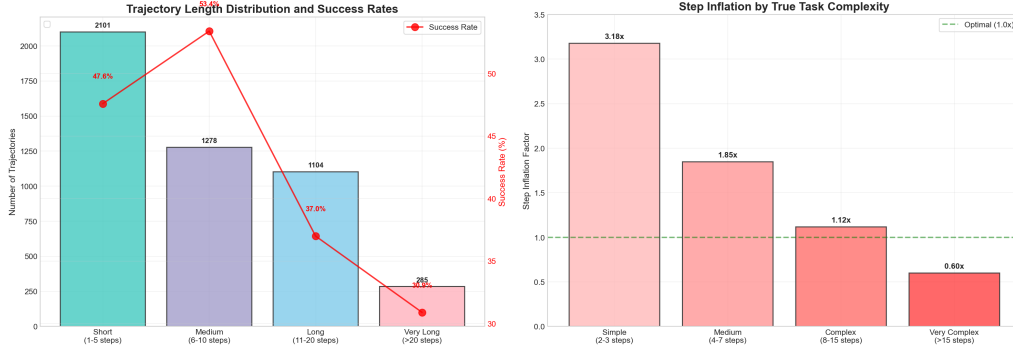


Figure 3: Trajectory efficiency and success patterns. (Left) Success rates follow an inverted-U curve with respect to trajectory length, peaking at medium trajectories (6–10 steps, 53.4%). (Right) Step inflation by task complexity shows a counter-intuitive trend: agents inflate simple tasks the most (3.18 \times), while complex tasks approach or even beat near-optimal paths due to survival bias.

Figure 3 (right) analyzes step inflation across complexity levels, defined as the ratio between an agent’s actual steps and the shortest successful path. Surprisingly, simple tasks show the largest inflation (3.18 \times), suggesting agents frequently take unnecessarily long detours even when optimal solutions are short. In contrast, complex and very complex tasks display lower inflation (1.12 \times and 0.60 \times), likely because only highly efficient trajectories survive, creating a selection bias toward near-optimal solutions. Together, these results show that both trajectory length and inflation provide complementary signals about structural efficiency.

4.4 Action Necessity Analysis

Overall, 76.7% of actions are labeled necessary by our LLM-based annotation system, showing that most behaviors directly contribute to task completion while a substantial 23.3% are exploratory or redundant. Framework-level differences are modest but notable: necessity rates range from 72.9% (Learn by Interact) to 82.0% (UI-TARS). Importantly, high necessity alone does not guarantee success—UI-TARS has the highest necessity rate but only a 38.7% success rate—indicating that efficient action selection must be paired with correct decision-making at key points.

Necessity varies systematically by action type. Type actions are most often essential (82.0%), followed by `Select` (79.7%) and `Click` (78.2%), while “other” actions are least necessary (70.2%). This suggests that data-entry operations nearly always advance task goals, whereas miscellaneous navigations are more prone to redundancy.

Task complexity and action confidence provide additional perspectives. As tasks grow more complex, the proportion of unnecessary actions rises steadily, from 15.8% in simple tasks to 31.1% in very complex tasks. Confidence strongly predicts necessity: high-confidence actions (>0.95) are necessary in 83.4% of cases, compared to only 61.7% for low-confidence actions (<0.85), with correlation $r = 0.67$, $p < 0.001$. Temporally, early actions are most critical (84.3% necessity in the first three steps), while later actions drop to 71.2%, reflecting the accumulation of detours and recovery attempts in longer trajectories.

Necessity is also dynamic rather than fixed. Figure 2 (right) shows that necessity rates increase with repeated task attempts, rising from 68% on first attempts to over 83% after ten. This learning curve demonstrates that agents can improve efficiency over time by reducing redundant actions. At the same time, anomalies such as one-step successes (13% of all successful trajectories) show that necessity can be misleading in edge cases, since some tasks can be completed with minimal effort due to action granularity or pre-filled states. Together, these results establish necessity as a measurable, learnable, and interpretable signal for evaluating agent efficiency, while also highlighting its limitations.

4.5 Behavioral Analysis

Consensus graphs provide a lens for understanding navigation strategies at both the path and action levels. At the trajectory scale, we observe three dominant styles of behavior. Direct navigation strate-

gies follow near-optimal paths with minimal detours and account for 43% of successful trajectories. Exploratory navigation strategies, seen in 31% of cases, involve backtracking and alternative path trials before arriving at the goal. The remaining 26% of trajectories exhibit hybrid approaches, where targeted actions are interleaved with selective exploration. This spectrum of strategies highlights that agents vary not only in whether they succeed but also in how they approach problem solving.

At the action scale, necessity analysis shows that 76.7% of actions are essential for task completion. Clear differences emerge across action types: Type actions are the most consistently necessary (82.0%), followed by Select (79.7%) and Click (78.2%), while miscellaneous actions are the least necessary (70.2%). Correlation with LLM confidence ($r = 0.68$) further indicates that redundant or exploratory actions are often accompanied by lower certainty. Temporally, earlier steps carry disproportionate importance, with 84.3% necessity in the first three actions versus only 71.2% after step ten. Longer trajectories therefore accumulate more redundancy and recovery attempts, making necessity a meaningful signal for distinguishing efficient from inefficient behaviors.

4.6 Cross-Agent Analysis

Moving from individual trajectories to collective behavior, we analyze cross-agent variation using entropy-based diversity metrics and clustering over trajectory distributions. Strikingly, 83.2% of tasks produce mixed outcomes across agents, underscoring that different frameworks bring complementary strengths rather than converging on a single strategy. Zeta Labs Jace.AI achieves its strongest performance on content creation tasks (66%), while IBM CUGA leads in structured updates (61%). Other agents reveal their own niches, reinforcing the view that no single framework dominates across task categories.

Consensus graphs further reveal shared structural backbones. Across models, 89% of successful trajectories begin with similar initial sequences, suggesting the presence of critical paths that anchor successful behavior. At the same time, 37% of failed trajectories terminate prematurely compared to their successful counterparts. This pattern implies that agents often recognize futility early, abandoning tasks after encountering traps rather than engaging in extended redundant exploration.

4.7 Key Findings

Our analyses yield several insights into web agent behavior. First, we observe a performance–efficiency trade-off: higher necessity rates signal more focused action sequences but do not always correlate with higher success, as illustrated by UI-TARS. Second, task complexity shows an inverted-U relationship with success, with medium-complexity tasks achieving the highest rates while both simple and very complex tasks prove more difficult. Third, frameworks exhibit clear complementarity, with different agents excelling in distinct categories and the majority of tasks showing heterogeneous outcomes across models. Finally, consensus graphs capture behavioral phenomena that single-path metrics overlook, such as the existence of shared critical paths and the tendency of failed trajectories to terminate early. Together, these findings demonstrate that WebGraphEval not only benchmarks performance but also surfaces the structural and behavioral dynamics that define agent strengths and weaknesses.

5 Discussion

In this work, we introduced **WebGraphEval**, a framework that extends evaluation beyond binary success rates by providing a structured, interpretable, and multi-dimensional view of web agent behavior. The results demonstrate that WebGraphEval captures not only whether agents succeed, but also how they navigate, where inefficiencies occur, and which strategies are shared or divergent across frameworks. This enables a richer understanding of strengths, weaknesses, and complementarities in current web navigation systems.

Despite these contributions, several limitations remain. First, the reliability of consensus graphs depends on the availability of diverse trajectories. Tasks with few attempts or limited agent coverage yield less stable structural insights. Second, the current state and action canonicalization is implemented through heuristics and LLM-based prompts. While effective in many cases, this approach may struggle with the breadth of real-world interfaces and actions. Third, contextual completeness is

constrained by the dataset: many trajectories lack full screenshots or auxiliary information, which limits the fidelity of environment reconstruction.

Looking forward, we see three main avenues for improvement. (1) **Reducing data dependence**: Few-shot graph construction and transfer learning across related tasks can extend the framework’s applicability to sparse or novel domains. (2) **Improving abstraction**: Replacing heuristic canonicalization with learned, semantically informed models can yield more robust state and action representations. (3) **Closing the loop**: Rather than being purely diagnostic, consensus graphs could inform online decision-making—either by guiding agent exploration during inference or by serving as a structured reward signal in reinforcement learning.

Acknowledgments

We thank Cookie (Yaoyao’s dog) and Lucas (Yaoyao’s cat) for their comforting presence during this work.

References

- [1] Shixiang Zhou, Biao Deng, et al. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2402.07762*, 2024.
- [2] Biao Deng, Qinglong Zhang, et al. Mind2web: Towards a generalist agent for the web. *arXiv preprint arXiv:2306.06070*, 2023.
- [3] Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, et al. A survey on llm-as-a-judge. *arXiv preprint arXiv:2411.15594*, 2024.
- [4] Shinn et al. Yao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- [5] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757, 2022.
- [6] 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.
- [7] Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan, and Dong Yu. Webvoyager: Building an end-to-end web agent with large multimodal models. *arXiv preprint arXiv:2401.13919*, 2024.
- [8] Boyu Gou, Zanming Huang, Yuting Ning, Yu Gu, Michael Lin, Weijian Qi, Andrei Kopanav, Botao Yu, Bernal Jiménez Gutiérrez, Yiheng Shu, et al. Mind2web 2: Evaluating agentic search with agent-as-a-judge. *arXiv preprint arXiv:2506.21506*, 2025.
- [9] Chang Ma, Junlei Zhang, Zhihao Zhu, Cheng Yang, Yujiu Yang, Yaohui Jin, Zhenzhong Lan, Lingpeng Kong, and Junxian He. Agentboard: An analytical evaluation board of multi-turn llm agents. *arXiv preprint arXiv:2401.13178*, 2024. NeurIPS 2024 (Oral).
- [10] Ido Levy, Ben Wiesel, Sami Marreed, Alon Oved, Avi Yaeli, and Segev Shlomov. Stwebagentbench: A benchmark for evaluating safety and trustworthiness in web agents. *arXiv preprint arXiv:2410.06703*, 2024.
- [11] Lawrence Jang, Yinheng Li, Charles Ding, Justin Lin, Paul Pu Liang, Dan Zhao, Rogerio Bonatti, and Kazuhito Koishida. Videowebarena: Evaluating long context multimodal agents with video understanding web tasks. 2024. Manuscript.
- [12] 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. *arXiv preprint arXiv:2308.03688*, 2023.

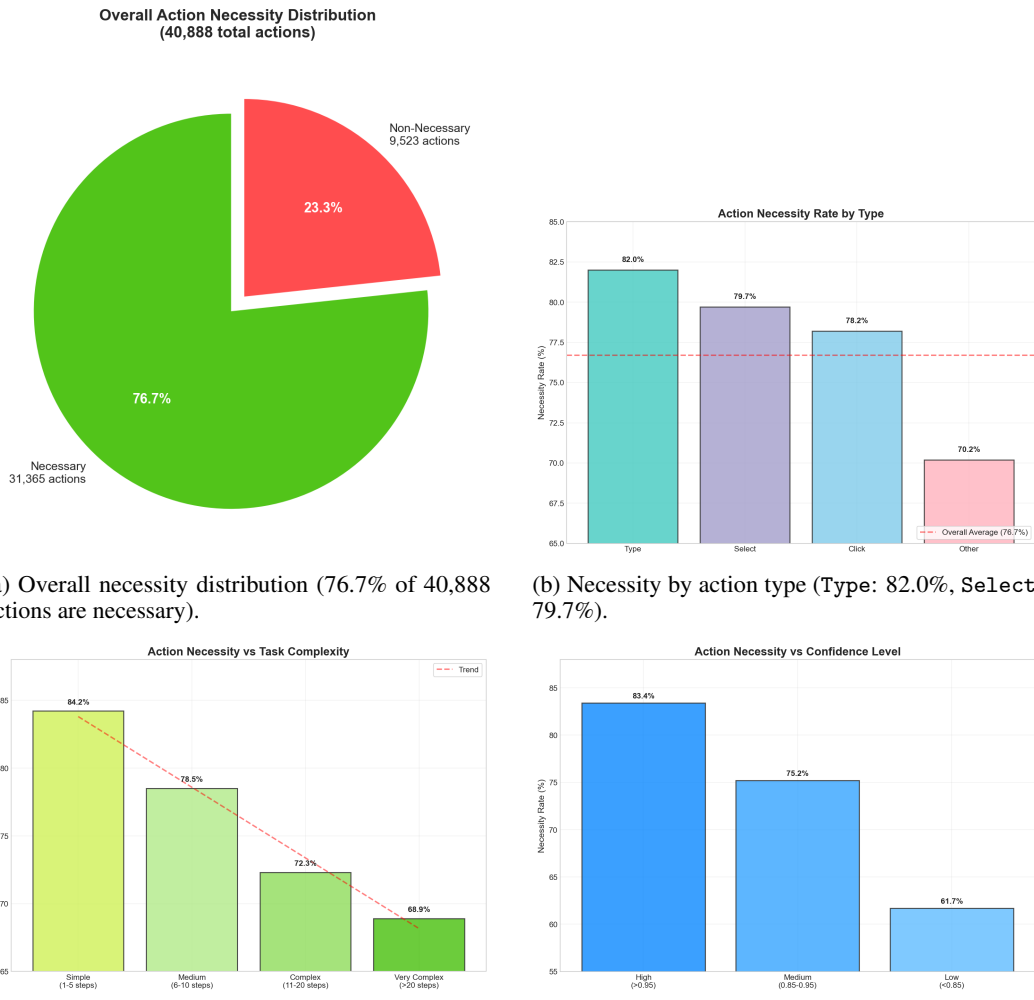
- [13] Seonghyeon Ye, Doyoung Kim, Sungdong Kim, Hyeonbin Hwang, Seungone Kim, Yongrae Jo, James Thorne, Juho Kim, and Minjoon Seo. Flask: Fine-grained language model evaluation based on alignment skill sets. *arXiv preprint arXiv:2405.12343*, 2024.
- [14] Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. *arXiv preprint arXiv:2010.03768*, 2020.
- [15] Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, Dasha Metropolitan, Robert Osazuwa Ness, and Jonathan Larson. From local to global: A graph rag approach to query-focused summarization. *arXiv preprint arXiv:2404.16130*, 2024.
- [16] Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao Huang. Lightrag: Simple and fast retrieval-augmented generation. *arXiv preprint arXiv:2410.05779*, 2024.
- [17] Abdul Haque, Ahmad Din, Muhammad Babar, Ali Abbas, Insaf Ullah, et al. Graphrag-causal: A novel graph-augmented framework for causal reasoning and annotation in news. *arXiv preprint arXiv:2506.11600*, 2025.
- [18] Haoran Luo, Guanting Chen, Qika Lin, Yikai Guo, Fangzhi Xu, Zemin Kuang, Meina Song, Xiaobao Wu, Yifan Zhu, Luu Anh Tuan, et al. Graph-r1: Towards agentic graphrag framework via end-to-end reinforcement learning. *arXiv preprint arXiv:2507.21892*, 2025.
- [19] Apurva Gandhi and Graham Neubig. Go-browse: Training web agents with structured exploration, 2025.
- [20] Jace AI Team. Jace ai: Email assistant that understands your voice. Online, 2025. Available at <https://jace.ai>.
- [21] IBM Research. Towards enterprise-ready computer using generalist agent. *arXiv preprint arXiv:2503.01861*, 2024. <https://arxiv.org/abs/2503.01861>.
- [22] Hongjin Su, Ruoxi Sun, Jinsung Yoon, Pengcheng Yin, Tao Yu, and Sercan Ö. Arik. Learn-by-interact: A data-centric framework for self-adaptive agents in realistic environments. *arXiv preprint arXiv:2501.10893*, 2025. <https://arxiv.org/abs/2501.10893>.
- [23] Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, et al. Ui-tars: Pioneering automated gui interaction with native agents. *arXiv preprint arXiv:2501.12326*, 2025. <https://arxiv.org/abs/2501.12326>.
- [24] OpenAI. Computer-using agent (cua). Online, 2025. Available at <https://openai.com/index/computer-using-agent/>.
- [25] Robin Gales and contributors. Browser-use: Make websites accessible for ai agents. <https://github.com/browser-use/browser-use>, 2025. GitHub repository.

Appendix

A Supplementary Figures

This appendix presents additional analyses and visualizations that complement the main text. To avoid redundancy, we organize related figures into compact multi-panel layouts.

A.1 Action Necessity Analysis



(a) Overall necessity distribution (76.7% of 40,888 actions are necessary).

(b) Necessity by action type (Type: 82.0%, Select: 79.7%).

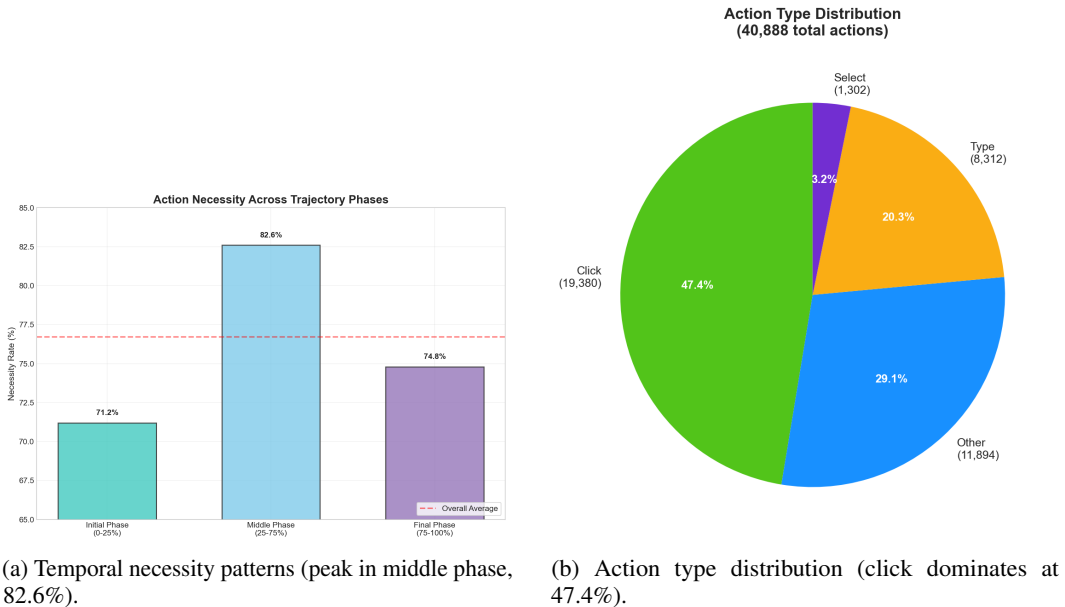
(c) Necessity vs. task complexity (drops from 84.2% to 68.9%).

(d) Necessity by action confidence ($r = 0.67$, $p < 0.001$).

Figure 4: Core necessity metrics across four dimensions: overall, action type, complexity, and confidence.

A.2 Framework Performance

A.3 Dataset Characteristics



(a) Temporal necessity patterns (peak in middle phase, 82.6%).

(b) Action type distribution (click dominates at 47.4%).

Figure 5: Supplementary behavioral patterns: (a) necessity varies across trajectory phases, with peak efficiency in the middle; (b) action type distribution reveals clicks dominate overall interactions.

Table 3: Overall Graph Dataset Statistics

Metric	Value
Total Graph Nodes	40,431
Total Graph Edges	45,656
Average Nodes per Task	49.79 ± 21.83
Average Edges per Task	56.23 ± 22.86
Average Steps per Trajectory	8.58

Table 4: Framework-level performance with mean ± standard deviation over 3 evaluation runs.

Framework	Success Rate (mean ± std)
Zeta Labs Jace.AI [20]	64.86% ± 0.43
IBM CUGA [21]	59.70% ± 0.14
Learn by Interact [22]	53.74% ± 0.43
UI-TARS [23]	38.88% ± 0.35
OpenAI-CUA [24]	28.55% ± 0.07
BrowserUse [25]	27.66% ± 0.15

Table 5: Comparison of graph-based approaches across domains. Prior work applies graphs for knowledge retrieval or causal reasoning, while WebGraphEval uses graphs to aggregate trajectories for evaluation.

Method	Graph Type	Node Type	Edge Type	Multi-hop	Method	Domain
GraphRAG [15]	Knowledge Graph	Entities	Semantic	✓	Traversal	Document QA
LightRAG [16]	Sparse Graph	Entities	Semantic	✓	Traversal	Real-time QA
Causal GraphRAG [17]	Causal Graph	Events	Causal	✓	Path-finding	News Analysis
Graph-R1 [18]	Hypergraph	Entities	Hyperedges	✓	Interaction	Complex QA
WebGraphEval (Ours)	Directed Graph	Canonical Actions	Transitions	✓	Trajectory Aggregation	Web Agent Evaluation

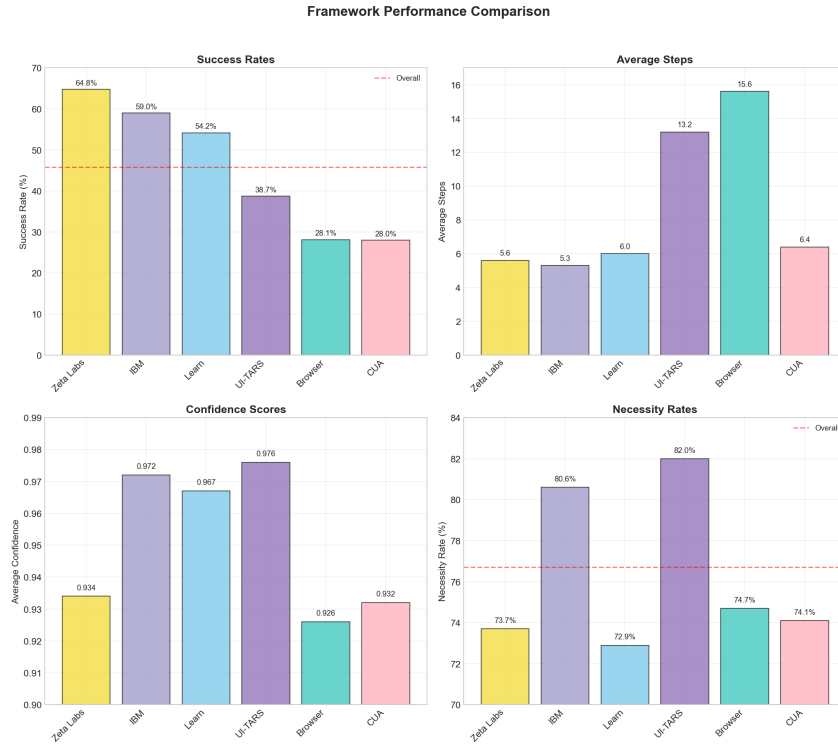


Figure 6: Comprehensive framework performance matrix comparing six web agent frameworks across success rate, average steps, confidence scores, and necessity rates.

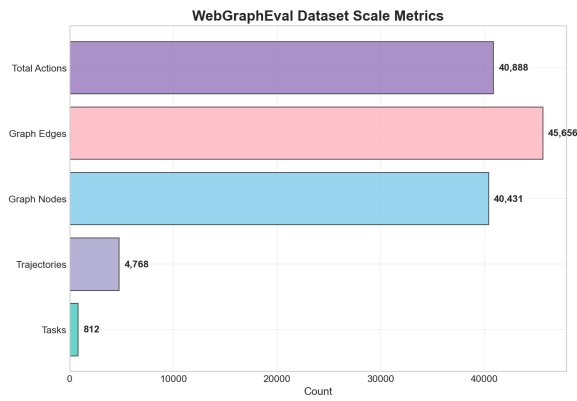


Figure 7: Dataset scale: 812 tasks, 4,768 trajectories, 40,431 nodes, 45,656 edges, and 40,888 actions.

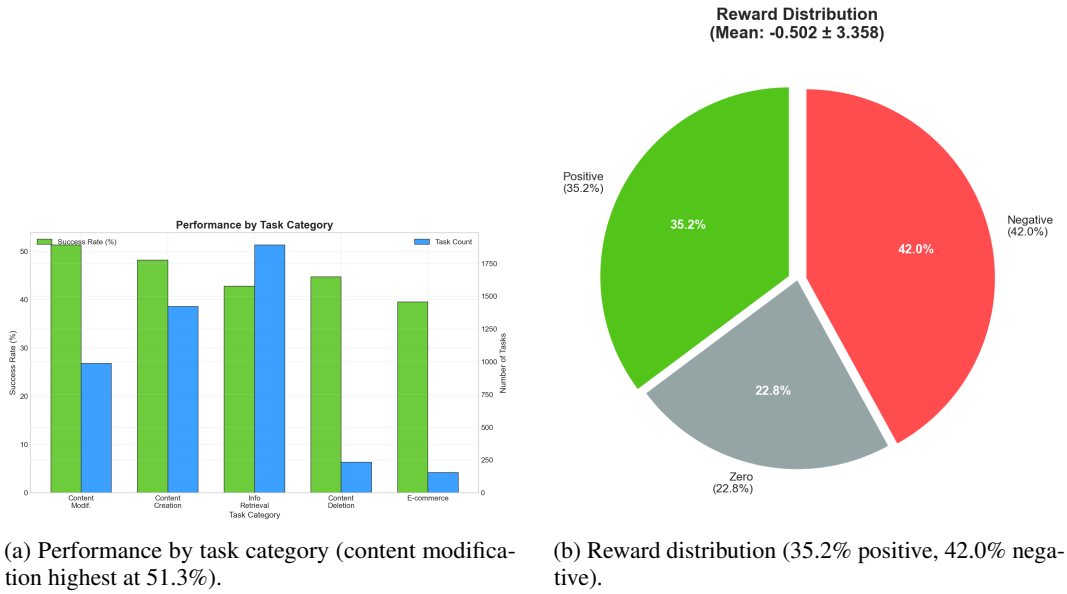


Figure 8: Task-level analysis by category and reward distribution.

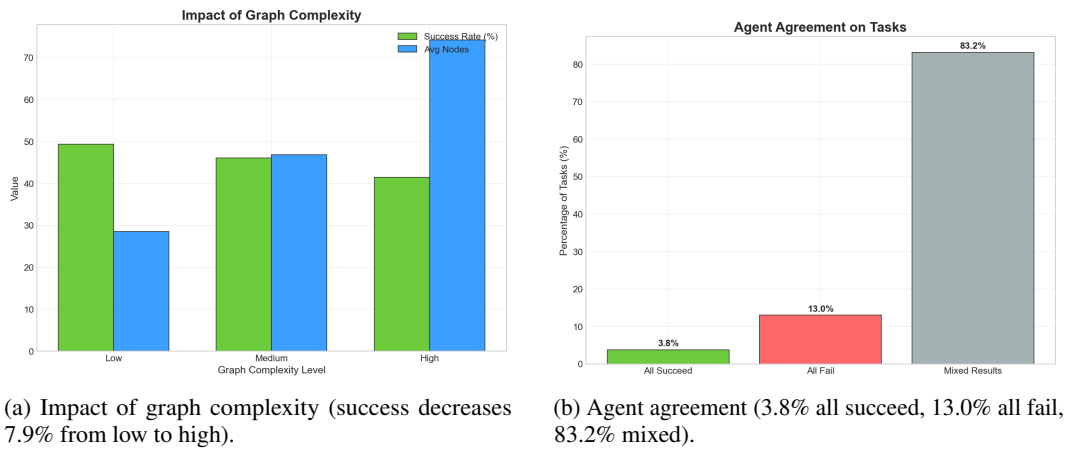


Figure 9: Complexity and agent agreement patterns.

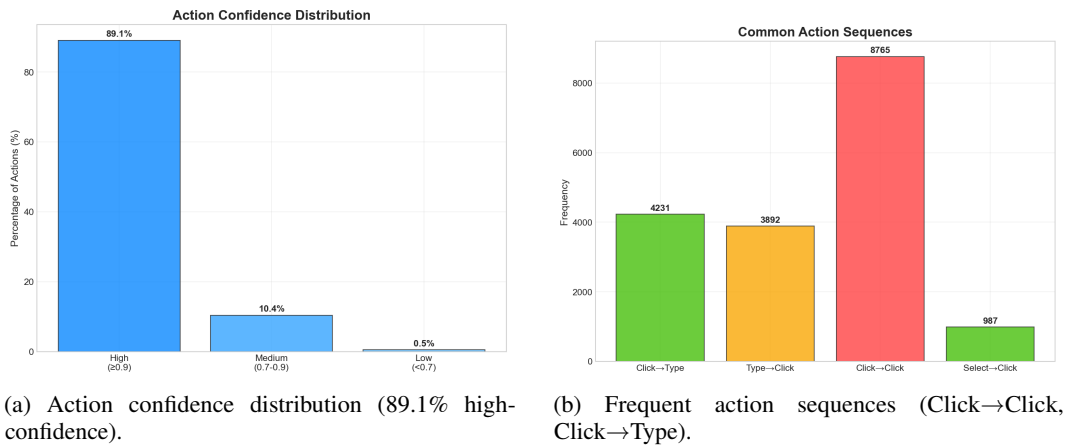


Figure 10: Confidence and sequence-level patterns.

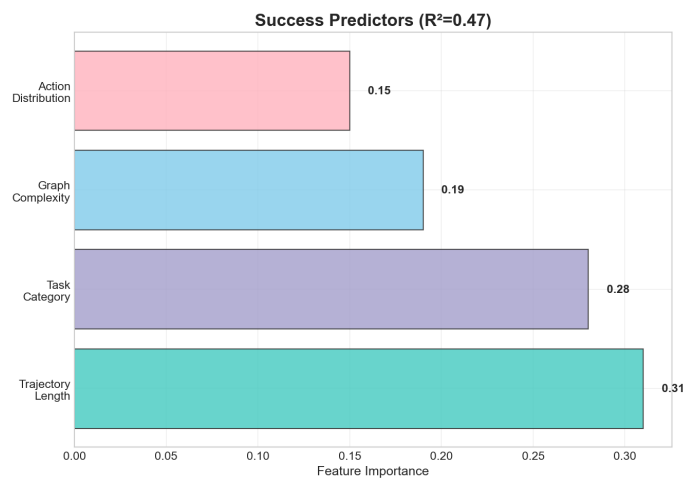


Figure 11: Success predictors: trajectory length (0.31), task category (0.28), graph complexity (0.19), action distribution (0.15). Overall R²=0.47.

Trajectory Success Evaluation

[System Prompt]

You are an expert evaluator for web agent trajectories. Your task is to determine if a web agent successfully completed a given task based on its actions and final response.

EVALUATION CRITERIA:

1. Analyze the agent's actions, URLs visited, and final message
2. Check if the agent fulfilled the task intent
3. If reference answers are provided, verify they appear in the trajectory
4. Consider both the journey (actions/URLs) and destination (final message)
5. Be strict but fair in your assessment

RESPONSE FORMAT:

Respond with exactly "SUCCESS" or "FAILURE" on the first line, followed by a brief explanation on the next line.

[User Prompt]

Task Intent: {intent}

Reference Answers:

- Exact match expected: "{exact_match}"
- Must include all of: "{must_include_items}"
- Fuzzy match acceptable: "{fuzzy_match_items}"

Final Message:

"{agent final reply}"

Action Sequence:

{index}. {action} (at {url})

...

EVALUATION CRITERIA (with / without Reference Answer):

1. Does the final message contain the correct answer that matches the reference?
2. If the answer is not in the final message, check if it appears in the actions or URLs
3. The reference answer **MUST** be found somewhere in the trajectory for success

IMPORTANT: The provided reference answers are the ground truth. Success requires finding these specific answers.

Respond with: SUCCESS or FAILURE, followed by a brief explanation.

Figure 12: Trajectory Success Evaluation prompt

System Prompt: Action Conversion and Necessity Annotation

You are an expert at converting natural language web action descriptions into standardized function calls.

Available Functions:

- `click(text: string, element?: string)` – Click on an element
- `type(text: string, element?: string)` – Type text into an input field
- `scroll(direction: "up"|"down", amount?: number)` – Scroll the page
- `select(value: string, element?: string)` – Select from dropdown
- `hover(text: string, element?: string)` – Hover over element
- `wait(seconds: number)` – Wait for specified time
- `goto(url: string)` – Navigate to URL
- `back()` – Go back in browser history
- `refresh()` – Refresh the page

Instructions:

1. Analyze the input action description carefully
2. Extract the key intent and parameters
3. Map to the most appropriate function from the available list
4. Use named parameters format: `functionName(param1="value1", param2="value2")`
5. If multiple interpretations are possible, choose the most likely one
6. Maintain high confidence for clear matches, lower for ambiguous cases
7. If no function matches well, return confidence < 0.5
8. ALSO decide whether the action is necessary for accomplishing the task (boolean necessary)
9. If the action depends on earlier steps, include a nested pre field: `{id:"<step id>", pre:{...}}`
10. If the input describes multiple discrete actions, split them and output a JSON array

Output Format (JSON):

```
{
  "functionName": "click",
  "parameters": ["Submit", "button"],
  "namedParameters": {"text": "Submit", "element": "button"},
  "confidence": 0.95,
  "necessary": true,
  "pre": { "id": "step 2" },
  "reasoning": "Clear click action on a button element"
}
```

Important:

- Always include both "parameters" (array) and "namedParameters" (object)
- Always include "necessary"
- Always include "pre" (null if no dependency)

Always output valid JSON. Be concise but accurate.

Figure 13: System prompt for action conversion and necessity annotation.

User Prompt: Action Conversion and Necessity Annotation

You are given a natural language description that may contain multiple action sentences. First split the text into individual sentences by period (.), exclamation (!), or question mark (?). For each non-empty sentence, produce a standardized function call. Output a JSON array of conversion objects, one per sentence.

Description:

“{action}”

Task Context: {task_description}

Previous Steps:

step 0 {previous_action_0}

step 1 {previous_action_1}

• ...

Figure 14: User prompt for action conversion and necessity annotation.