

Odysseys: Benchmarking Web Agents on Realistic Long Horizon Tasks

Lawrence Jang*, Jing Yu Koh*, Daniel Fried & Ruslan Salakhutdinov
Carnegie Mellon University
{ljang, jingyuk, dfried, rsalakhu}@cs.cmu.edu

Abstract

Existing web agent benchmarks have largely converged on short, single-site tasks that frontier models are approaching saturation on. However, real-world web use consists of long-horizon, multi-site workflows. Common web navigation tasks such as comparing products across different domains, planning trips across multiple services, or summarizing information from multiple search queries, require sustained context, cross-site reasoning, and coherent planning over potentially hours of browsing. To capture and evaluate such behaviors, we introduce Odysseys: a benchmark of 120 long-horizon web tasks derived from real world browsing sessions, evaluated on the live Internet. We find that binary pass/fail evaluation is inadequate for long-horizon settings, and introduce a rubric-based evaluation, annotating each Odysseys task with an average of 5.8 graded rubrics. We demonstrate that this yields higher agreement with humans, and provides a more fine-grained signal than commonly used trajectory-level LLM-as-a-judge evaluation metrics. We test several leading frontier models, and find that the strongest models achieve a success rate of 53%, which leaves substantial headroom for future improvements. Odysseys isolates the critical evaluation of long-horizon proficiency in open-web environments, providing a realistic benchmark to measure progress towards computer-use agents that can potentially productively operate for hours.

1 Introduction

Recent large language models can now function as computer-use agents. They browse websites, interpret screenshots, click interface elements, and execute multi-step instructions in human-facing software. However, current benchmarks largely evaluate these capabilities in short, tightly scoped episodes, leaving a key regime underexplored: long-horizon web workflows that unfold across many pages, tabs, and domains. This limitation matters because real-world web use is rarely confined to a single page or a single domain. Many common tasks require extended interaction over multiple sites, such as comparing products across retailers, planning travel across booking platforms, or gathering information from search results and synthesizing it into a downstream deliverable. Solving such tasks requires more than local grounding or short-term action selection. Agents must maintain context over long horizons, reason across heterogeneous websites, decompose open-ended goals into subproblems, and decide when to stop exploring and produce an output.

Towards this goal, we introduce Odysseys, a benchmark of 120 long-horizon web tasks derived from real-world browsing behavior and evaluated on the live Internet. Each task require agents to execute multi-step workflows across multiple websites. The tasks cover realistic activities, and are grounded in annotated human browsing journeys rather than synthetic templates. We also identify the inadequacy of trajectory-level LLM-as-a-judge evaluation metrics (He et al., 2024; Xue et al., 2025), commonly used in computer-use benchmarks. These metrics become increasingly noisy as trajectories grow longer and more complex. To address this, we propose a rubric-based evaluation scheme in which each

*Equal contribution

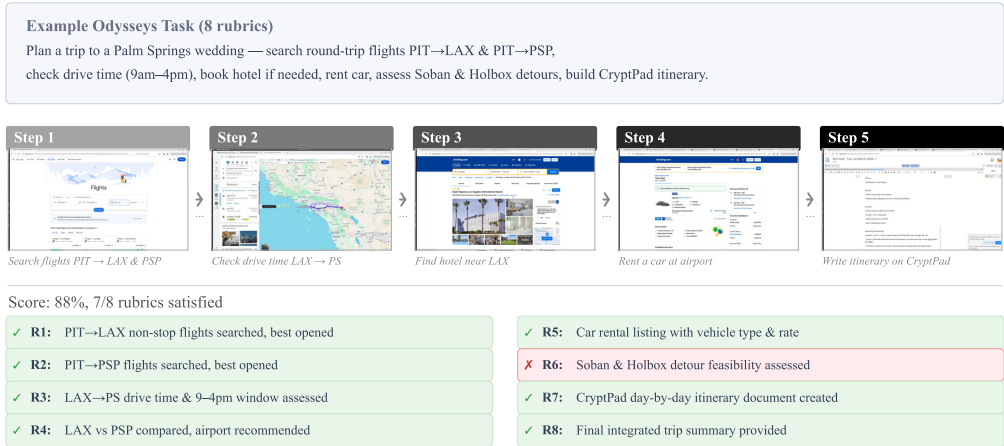


Figure 1: Odysseys is a long-horizon web agent benchmark with 120 tasks based on real user browsing data. Example task (with simplified task description), trajectory and rubric-based evaluation visualized above.

task is decomposed into a set of verifiable requirements. This provides a more informative measure of partial completion and yields stronger agreement with human judgments than trajectory-level grading.

We evaluate leading frontier and open-weight computer-use agents on Odysseys. The strongest model we tested (Opus 4.6) achieves 53% perfect task success, indicating substantial headroom for progress. Our analysis shows that current agents struggle with core aspects of long-horizon web navigation, including sustained planning, cross-site coordination, and balancing information gathering against completing user requested deliverables. At the same time, we observe several compelling strategies emerge, suggesting that frontier agents have developed flexible behaviors for operating in open-ended computer environments.

We conduct an analysis of model performance, test-time scaling with more steps, and do a qualitative analysis and identify common failure modes, highlighting key challenges for future work on building robust long-horizon computer-use agents. Our benchmark is publicly released and maintained at `removed_for_review`.

2 Related Work

Computer-Use Agents Benchmarks for computer-use agents based on large language models (LLMs) began with synthetic environments such as MiniWoB (Shi et al., 2017), MiniWoB++ (Liu et al., 2018) and WebShop (Yao et al., 2022). More recently, several realistic benchmarks, such as OSWorld (Xie et al., 2024b), Windows Agent Arena (Bonatti et al., 2025), and MacOSWorld (Yang et al., 2025), measure the performance of LLM-based agents on OS-level utilities and office suites. AndroidWorld (Rawles et al., 2025) and WorkArena (Drouin et al., 2024) extend coverage to mobile and enterprise domains respectively, while AgentBench (Liu et al., 2024) spans web, database, and OS environments.

In particular, for web browsing, several benchmarks have been proposed to evaluate GUI-based navigation and task completion, requiring agents to interact through screenshots, clicking, and typing. Mind2Web (Deng et al., 2023) introduced the evaluation of LLM agents in web-based tasks using static webpages. WebArena (Zhou et al., 2024) and VisualWebArena (Koh et al., 2024) created dynamic synthetic web environments with execution-based verification that is precise but hard to scale. Online-Mind2Web (Xue et al., 2025) and Web-Voyager (He et al., 2024) measure progress with LLM-as-a-judge over full trajectories on real sites in the live internet, which scales but introduces grading noise. Navi-Bench (Yutori, 2025) alleviates some of these issues by introducing a dynamic task configuration that renders queries and success criteria based on the current date and state of live websites. Several

static web grounding benchmarks, such as ScreenSpot (Cheng et al., 2024; Li et al., 2025), OmniACT (Kapoor et al., 2024), and OSWorld-G (Xie et al., 2025), also provided evaluation measures for visual grounding specific metrics. At longer research horizons, GAIA (Mialon et al., 2024) and BrowseComp (Wei et al., 2025) evaluate multi-hop information synthesis from the open web, targeting the deep research paradigm, but shift away from GUI-level interaction toward API-based or text-based access, evaluating final answers rather than the browsing process itself. Critically, all of these web benchmarks target short-answer or single-site tasks (which are saturated by frontier models), or do not truly test the full multimodal and navigation capabilities of a computer-use agent.

Long-Horizon LLM Agents As LLMs have increased their general capabilities to reason and act as autonomous agents, their ability to expand into long-horizon tasks has become a central focus. METR (Kwa et al., 2026) introduces a task-completion time horizon as a measure of long-horizon capability, showing that the effective task length of frontier agents (measured by human completion time) has grown exponentially, roughly doubling every seven months from 2019–2025 with signs of further acceleration. Despite this progress, current agents remain limited to tasks on the order of hours and still fall short of reliably completing longer, multi-hour or day-scale workflows. Early work on multi-step reasoning of LLMs, such as HotpotQA (Yang et al., 2018) and MuSiQue (Trivedi et al., 2022), studied multi-hop reasoning over multiple documents but mainly operated in text, static settings without sequential decision-making or environment interaction. Several long-horizon agentic benchmarks have also been released: τ -bench (Yao et al., 2025) measures consistency in multi-turn tool calling agents; Vending-Bench (Backlund & Petersson, 2025) tests LLM agents on their ability to run a vending machine over a long time period; TheAgentCompany (Xu et al., 2025) evaluates multi-platform workplace tasks in a simulated software company; and TravelPlanner (Xie et al., 2024a) evaluates the collapse of LLMs in large-scale multi-constraint planning. OdysseyBench (Wang et al., 2025) extends this line to long-horizon office workflows across multiple applications.

3 The Odysseys Dataset

We introduce the Odysseys dataset, consisting of 120 long-horizon, multi-site web tasks designed to evaluate web agents on realistic browsing workflows. Each task begins from a Google search page and requires the agent to navigate across multiple websites to complete complex, multi-step objectives such as comparing products across e-commerce platforms, planning travel itineraries, or setting up video playlists and watching lectures. We show an overview of a few example tasks in Tab. 2.

3.1 Collection Process

We recruit 248 participants based in the US and UK through the Prolific¹ crowdsourcing platform and provide them with a desktop application that reads their Chrome browsing history and provides the user with an interface to annotate their own history into web agent tasks. The application segments browsing activity into singular tasks with corresponding URLs using Chrome’s Journey algorithm (Google, 2022). For each journey, participants annotate: (1) a *key URL* representing the success state, (2) an *automation preference* indicating whether they would want the task automated, (3) a *task label* describing the activity as they would prompt an AI tool, and (4) a *feasibility* judgment that marks whether the task would be possible to complete. The collection yields 2,380 labeled journeys that span various domains, from comparison shopping and travel planning to media consumption and information research. A visualization of our annotation interface is provided in Appendix A.2.

3.2 From Journeys to Odysseys

Raw journey labels are noisy, as participants may submit labels unrelated to the URLs, underspecify tasks, or provide vague descriptions. We address this through a two-stage

¹<https://prolific.com>

Level	Task Summary	Example Rubric Criteria
Easy	Find a roasted Brussels sprouts recipe with Parmesan and balsamic vinegar; check Le Creuset for a 5.5 qt Dutch oven in light green; find a YouTube video on dishwasher air gap installation.	<ul style="list-style-type: none"> – Recipe page shows both Parmesan and balsamic vinegar with title, oven temp, and cook time – Le Creuset product page confirms whether 5.5 qt is available in light green
	Build a weather risk snapshot: get Baltimore current conditions, Syracuse 7-day low from Wunderground, NWS forecast and alerts for Rittman, OH, and predicted snowfall from NWS Mount Holly.	<ul style="list-style-type: none"> – Report the lowest forecasted temperature in Syracuse over the next 7 days from Wunderground – Report predicted snowfall from the NWS Mount Holly winter forecast graphic
Medium	Plan a trip to Japan as a Canadian in the U.S.: find the passport renewal process, fees, and processing times; check visa requirements; review travel advisories; compare travel insurance from PolicyAdvisor and Kanetix.	<ul style="list-style-type: none"> – Correctly identify the official Canadian passport renewal process for citizens abroad – Determine whether a Canadian passport holder needs a tourist visa for Japan
	Evaluate a Tesla Model 3 lease in Los Angeles: find current lease terms and price, search Craigslist for trailer alternatives less than \$10K, find a comparable Zillow rental, and present a cost comparison.	<ul style="list-style-type: none"> – Summarize current LA Tesla Model 3 lease offer with payment, term, and due-at-signing – Identify ≥ 3 Craigslist trailer listings under \$10K suitable for living
Hard	Plan a trip to all 30 MLB stadiums: select one real game date per stadium from the official schedule, prioritize star players, plan travel logistics, find accommodations and local food via Yelp, and compile in CryptPad.	<ul style="list-style-type: none"> – Itinerary includes all 30 stadiums with one official game date and matchup each – Complete visit sequence with cheaper travel mode identified between all stops
	Find the top 10 banana bread recipes by popularity, compare ingredients and methods, select 3 most unique, find associated YouTube videos, and synthesize a combined best recipe explaining design choices.	<ul style="list-style-type: none"> – Ingredients and cooking methods extracted and compared across all 10 recipes – Combined recipe created drawing from strengths of top 3, with complete instructions

Table 1: Representative tasks from the Odysseys benchmark spanning three difficulty levels and diverse real-world domains. Task descriptions are condensed for readability to capture key actions (the full descriptions are available in Appendix A.6). Each task requires the agent to navigate multiple websites, extract specific information, and satisfy independently verifiable rubric criteria.

refinement process combining LLM screening with manual review by the authors (see Appendix A.3 for details). After filtering for label accuracy, feasibility, login requirements, and overall quality, 696 high quality journeys (29.2%) remained.

The resulting journeys represent real browsing behavior but are generally short subtasks due to Chrome’s journey segmentation. To compose them into realistic long-horizon tasks, we cluster related journeys by embedding similarity and use GPT 5.4 to arrange subsets into coherent multi-step workflows (see Appendix A.4 for full details of the clustering, composition pipeline, and prompt). For each composed task, the LLM generates a natural-language prompt, step plan, rubric with verification procedures, and a self-assessed coherence score. Each source journey is used at most 3 times across all Odysseys to prevent over-representation, and the authors perform manually quality assurance by inspecting every long chained Odyssey and de-duplicating any tasks that appear too similar (see 6). We reject generated tasks that span fewer than 2 sites, have incomplete rubrics, or score below 2/5 on coherence. Difficulty is assigned by step count and domain spread: *easy* (≤ 5 steps, ≤ 3 domains), *medium* (6–8 steps or 4+ domains), *hard* (exceeding both). In addition to the 90 composed tasks, the authors of this paper contributed 30 manually authored long-horizon tasks grounded in real personal queries, following the same schema for a total of 120 tasks.

Level	Task Summary	Example Rubric Criteria
Easy	Find a roasted Brussels sprouts recipe with Parmesan and balsamic vinegar; check Le Creuset for a 5.5 qt Dutch oven in light green; find a YouTube video on dishwasher air gap installation.	<ul style="list-style-type: none"> – Recipe page shows both Parmesan and balsamic vinegar with title, oven temp, and cook time – Le Creuset product page confirms whether 5.5 qt is available in light green
	Build a weather risk snapshot: get Baltimore current conditions, Syracuse 7-day low from Wunderground, NWS forecast and alerts for Rittman, OH, and predicted snowfall from NWS Mount Holly.	<ul style="list-style-type: none"> – Report the lowest forecasted temperature in Syracuse over the next 7 days from Wunderground – Report predicted snowfall from the NWS Mount Holly winter forecast graphic
Medium	Plan a trip to Japan as a Canadian in the U.S.: find the passport renewal process, fees, and processing times; check visa requirements; review travel advisories; compare travel insurance from PolicyAdvisor and Kanetix.	<ul style="list-style-type: none"> – Correctly identify the official Canadian passport renewal process for citizens abroad – Determine whether a Canadian passport holder needs a tourist visa for Japan
	Evaluate a Tesla Model 3 lease in Los Angeles: find current lease terms and price, search Craigslist for trailer alternatives less than \$10K, find a comparable Zillow rental, and present a cost comparison.	<ul style="list-style-type: none"> – Summarize current LA Tesla Model 3 lease offer with payment, term, and due-at-signing – Identify ≥ 3 Craigslist trailer listings under \$10K suitable for living
Hard	Plan a trip to all 30 MLB stadiums: select one real game date per stadium from the official schedule, prioritize star players, plan travel logistics, find accommodations and local food via Yelp, and compile in CryptPad.	<ul style="list-style-type: none"> – Itinerary includes all 30 stadiums with one official game date and matchup each – Complete visit sequence with cheaper travel mode identified between all stops
	Find the top 10 banana bread recipes by popularity, compare ingredients and methods, select 3 most unique, find associated YouTube videos, and synthesize a combined best recipe explaining design choices.	<ul style="list-style-type: none"> – Ingredients and cooking methods extracted and compared across all 10 recipes – Combined recipe created drawing from strengths of top 3, with complete instructions

Table 2: Representative tasks from the Odysseys benchmark spanning three difficulty levels and diverse real-world domains. Task descriptions are condensed for readability to capture key actions (the full descriptions are available in Appendix A.6). Each task requires the agent to navigate multiple websites, extract specific information, and satisfy independently verifiable rubric criteria.

3.3 Rubrics

Prior work on computer-use agents typically leverages execution-based verification (Zhou et al., 2024; Koh et al., 2024; Xie et al., 2024b) which requires handcrafted task-specific rewards, or LLM-as-a-judge over web agent trajectories (He et al., 2024; Xue et al., 2025) where an LLM takes in the screenshots and actions of an agent execution trace, and outputs a score or a pass/fail metric. However, we find that these are unsuitable for Odysseys, due to the long-horizon nature of the tasks, as well as the difficulty of crafting rules-based rewards for each task. Inspired by prior work on rubrics and checklists for evaluating and training language models (Shao et al., 2025; Hashemi et al., 2024; Gunjal et al., 2026; Viswanathan et al., 2025), we investigate using rubric-based evaluation for Odysseys tasks.

For each task, we generate a structured set of rubrics for evaluating partial progress (examples in Fig. 1, Fig. 2, and Tab. 2). Rubrics are generated alongside the task during LLM composition. For each rubric item, the model produces a *requirement*, a single verifiable checkpoint), and a *verification* description of how a grader should determine whether it is accomplished. Tasks contain 3–12 rubric items, with an average of 5.8 rubrics per task. All rubric items are verified by an author for correctness and relevance to task completion, and we also run an LLM augmented with text web search to cross-check verification criteria against live websites.

Type	Model	O-M2W Judge (↑)	Rubrics (%) (↑)		SPL (%) (↑)	
			Averaged	Perfect	Averaged	Perfect
API	Opus 4.6	43.8	82.9	52.9	1.36	0.98
	GPT-5.4	39.7	80.2	48.3	1.62	0.97
	Sonnet 4.6	35.5	76.3	43.8	1.13	0.75
	GPT-5.4-mini	36.5	67.4	27.0	1.59	0.61
Open Weights	Qwen-3.5-30B-A3B	16.8	44.7	13.3	0.69	0.23
	Qwen-3.5-9B	24.3	57.3	6.1	1.11	0.13
	UI-TARS-1.5-7B	4.3	38.0	2.5	0.61	0.05

Table 3: Model performance on Odysseys.

3.4 Final Review

After chaining, all task prompts are “CUAified”: Every task is rewritten from abstract information-gathering descriptions into natural, conversational requests styled as messages to a computer-use agent. Additionally, the LLM rewrites each prompt to include specific sites with concrete parameters (e.g., zip codes, price ranges, filter criteria), first-person motivation, natural step chaining without numbered lists, and visual browser-proof requirements such as keeping pages open in tabs, adding items to carts, or opening side-by-side comparisons. The rubric, steps, dependencies, and deliverable are regenerated in the same pass to align with the rewritten prompt. This ensures that each task reads as a realistic user request and that the evaluation criteria reference observable browser states rather than abstract information retrieval.

Every task, composed or manually authored, undergoes a two-pass review using a dedicated QA interface (Fig. 6). In the first pass, the authors verify prompt coherence, trace each component to its source journey, and ensure that rubric items are unambiguous and actionable. All prompts that reveal any potential personally identifiable information (PII) are also rewritten or removed. In the second pass, an LLM-augmented pass uses a web search tool to confirm site accessibility and step feasibility; a human adjudicates flagged issues. Tasks that fail in either pass are revised or removed. After post-processing, the Odysseys benchmark contains 120 tasks.

3.5 Dataset Statistics

The benchmark comprises 120 tasks spanning three difficulty levels: easy (45), medium (46), and hard (29). Task instructions average 265.8 words (median 264.5, range 76–387), reflecting the detailed, multi-step nature of real-world web workflows. Tasks are evaluated using a total of 699 rubric items, with an average of 5.8 rubrics per task (median 6, range 3–12). The dataset spans 21 top-level domains and 71 fine-grained categories from SimilarWeb, with each task touching an average of 2.0 categories (up to 3). The most represented domains are Travel and Tourism (34 tasks), Science and Education (33), Ecommerce and Shopping (25), and Computers, Electronics and Technology (25), followed by Health (17), Food and Drink (15), Arts and Entertainment (14), and Community and Society (14). This broad coverage ensures that the benchmark tests cross-site reasoning across diverse real-world contexts.

4 Experiments

We benchmark several frontier API-based models as well as open weights models on Odysseys. Each model is implemented with the recommended settings for computer-use, and we launch them in the same virtual Ubuntu environments from OSWorld (Xie et al., 2024b). The primary application that the models use is Google Chrome, although some models also occasionally choose to use other applications (such as LibreOffice for generating a report). We only benchmarked models that have general computer-use abilities, and do not test models that are only trained for web navigation, as many Odysseys tasks involve subtasks that are infeasible for web navigation only models (e.g., maintaining multiple tabs).

Metric	Cohen's κ (\uparrow)	F1 (\uparrow)	Accuracy (\uparrow)
O-M2W Judge	0.501	0.730	0.742
Rubrics (Averaged)	0.705	0.935	0.898
Rubrics (Perfect)	0.717	0.870	0.860

Table 4: Automated LLM-based metrics vs. human agreement at three granularities. Rubrics scores (Cohen's κ , F1, and accuracy against human labels) show substantial agreement, while the trajectory-level Online-M2W (holistic binary judgment) agreement is worse.

Task #16 (6 rubrics)
 I'm trying to get a realistic shortlist of the best knee surgeons in New York City because I may need ACL reconstruction or meniscus repair, and I want something I can actually look through myself afterward. Please start in Google and research orthopedic surgeons in NYC who are specifically known for knee ligament reconstruction, ACL surgery, and meniscus repair, then create a spreadsheet called Top ACL Surgeons NYC to keep everything organized. As you find strong candidates, open each surgeon's official hospital or practice profile page in its own tab so I can compare them side by side, and only keep surgeons whose actual profile page clearly says they perform ACL reconstruction, meniscus repair, knee ligament reconstruction, or very closely related sports knee procedures. For each surgeon you keep, put their full name, hospital or practice affiliation, specialty focus, a short note confirming where ACL reconstruction or meniscus repair is mentioned, and the direct profile link into the spreadsheet. Please keep going until there are exactly 10 verified NYC surgeons in the sheet, and make sure every person listed still has their real profile page open in a tab so I can inspect the pages and see the affiliations myself. Once the list is complete, look across the 10 entries and add a short summary of which hospitals, orthopedic groups, or medical centers show up most often, because I want to know which institutions seem to dominate this specialty in the city.

GPT-5.4 | Score: 100%, 93 steps

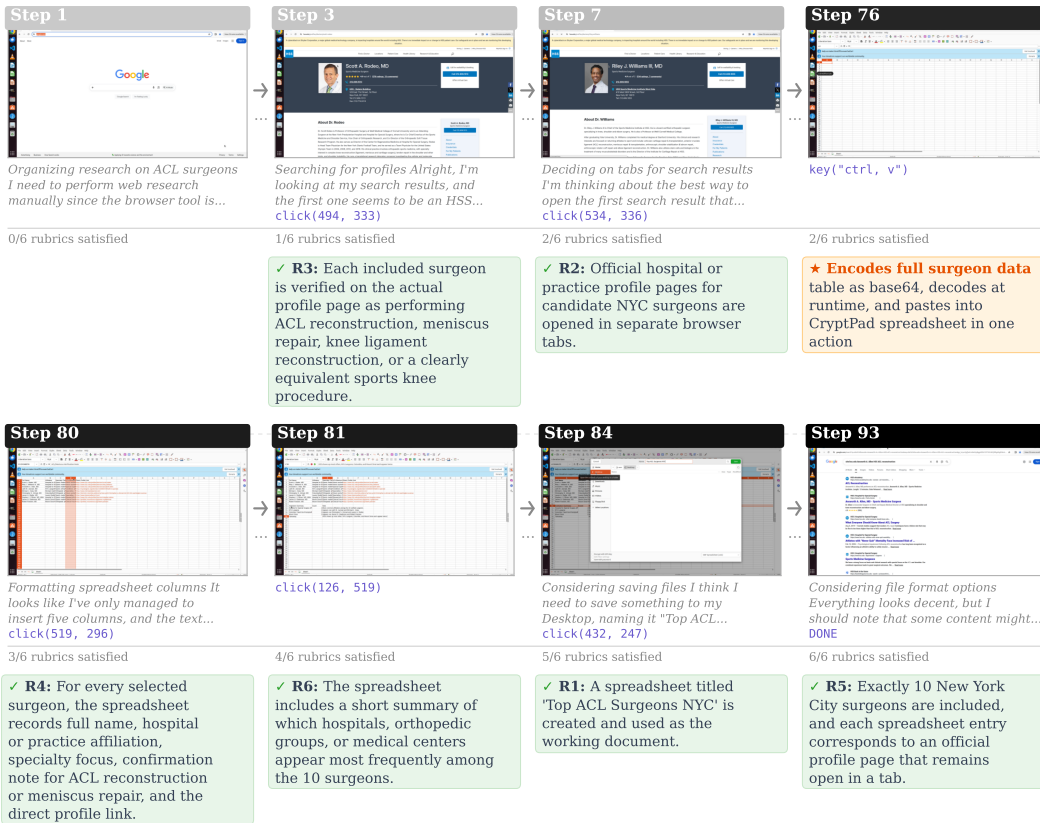


Figure 2: After researching 10 surgeon profiles, GPT-5.4 encoded the entire data table as a base64 string inside the Python action, decoded it at runtime with `pyperclip.copy(_text); pyautogui.hotkey('ctrl', 'v')`. This preserves exact tab and newline characters, and populates the spreadsheet in a single paste with the correct column mapping.

4.1 Evaluation

Rubrics Scores As discussed in Sec. 3.3, we compute rubrics scores based on the specific success criteria of each task. To ascertain the effectiveness of this metric, we measure agreement between automated LLM judges and human annotations at three levels of granularity, summarized in Tab. 4. At the finest level, *rubrics (average)* treats each task-rubric pair as an independent observation and computes agreement across all 699 rubric pairs; for the *rubrics (perfect)* metric, a task is scored as passing only if every one of its rubrics are satisfied, yielding a single binary label per task; and at the *trajectory* level, we run the LLM judge from Online-M2W (Xue et al., 2025), which issues a holistic pass/fail judgment based on the agent’s action trajectory. We collected human annotations of all 120 Opus 4.6 trajectories, where each annotator rated whether each rubric was satisfied during the course of the trajectory (see Appendix A.7 for a visualization of the annotation interface). Then, we compute human agreement of each metric by measuring Cohen’s κ coefficient, the F1 score, and the accuracy against the collected human labels. Across all metrics, the rubric-based evaluations substantially outperform the trajectory judge: rubric-average agreement achieves a Cohen’s κ of 0.705 and F1 of 0.935, compared to 0.501 and 0.730 for the trajectory judge. The cumulative task-level metric shows even stronger agreement ($\kappa = 0.717$, $F1 = 0.870$), demonstrating that the rubric-based approach remains reliable even under the stricter all-or-nothing criterion. Our proposed rubric-based evaluation (decomposing task completion into fine-grained, state-grounded rubrics) yields more reliable automated judgments than holistic trajectory-based assessment, especially for the long-horizon tasks in Odysseys. Additionally, rubrics also provide partial credit for weaker models, which may only achieve zero scores under the trajectory-level judge. This also introduces a meaningful reward signal for training weaker models with reinforcement learning, which would be a strong direction to explore in future work.

SPL (Success weighted by Path Length). Raw rubric scores measure what an agent accomplished but not *how efficiently* it did so. Two agents may both complete a task, yet one may take 30 steps while the other takes 100. To capture this dimension of Odysseys runs, we report SPL, adapted from the navigation metric of Anderson et al. (2018). For all task i in Odysseys, we compute

$$\frac{1}{N} \sum_{i=1}^N \frac{s_i}{n_i} \tag{1}$$

where s_i is the rubric score (either averaged or perfect) and n_i is the number of agent steps. The main difference in our metric is that unlike tasks in Vision-Language Navigation, we do not know the oracle number of steps required to complete a task, and therefore cannot normalize by this. Higher SPL indicates that an agent achieves strong outcomes in fewer steps, penalizing inefficient trajectories that reach the same result through unnecessary actions.

4.2 Results

We run all experiments under the same settings: we limit the number of steps to at most 100 steps², and we use the max reasoning effort possible (if any). We use the runner from OSWorld (Xie et al., 2024b), and start the models off with a Google Chrome window at the provided starting url for each task (if any). The models have access to the full Google Chrome interface, and are allowed to open tabs, save files, or do anything else supported by the virtual machine. We evaluate all models with the proposed rubric-based evaluation (Sec. 3.3). The results are summarized in Tab. 3. We observe that the best frontier models trained with computer use capabilities (Opus 4.6, GPT-5.4) achieve the strongest scores, with Opus 4.6 scoring slightly better across most metrics. However, GPT-5.4 (and GPT-5.4-mini) have higher SPL scores, indicating that they are more efficient at accomplishing Odysseys tasks. Open weight models trail starkly behind frontier models, with the Qwen-3.5-30B-A3B model (the largest we were able to run locally) performing the best amongst open weights computer use agents.

²For 100 steps, Opus 4.6 takes approximately 30 mins and \$2.5 to run per task.

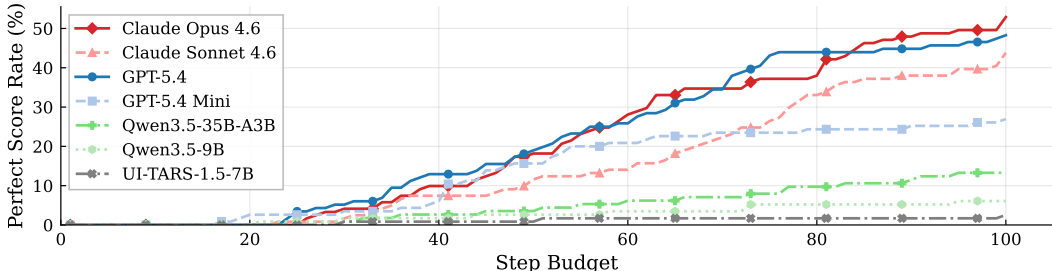


Figure 3: Perfect rubric rate as a function of step budget. Each curve shows the mean rubrics score achieved within a given number of interaction steps.

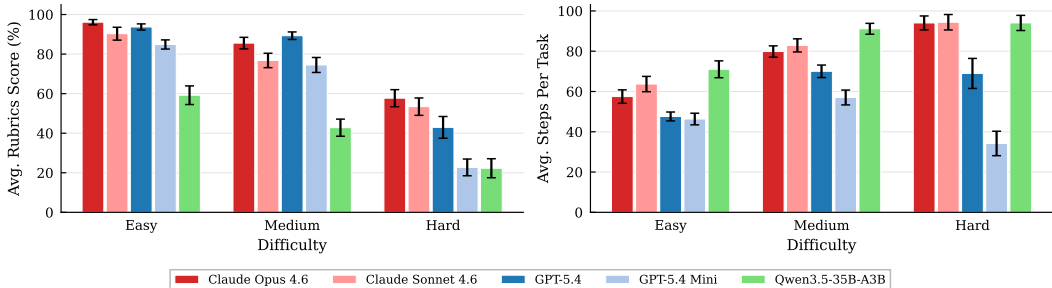


Figure 4: Average rubric scores (left) and # of steps taken per task (right) broken down by difficulty.

4.3 Analysis

In this section, we perform a careful analysis of the results and highlight several interesting findings, as well as possible directions for future work.

Scaling with step budget Fig. 3 shows how task completion improves as models are allowed more interaction steps. All models exhibit a broadly sigmoidal scaling curve: perfect score rates remain near zero for the first ~15 steps, reflecting the minimum number of interactions needed to complete the simplest Odysseys tasks. Rates then rise steadily through the 20–70 step range as progressively harder tasks are fully solved, before the rate of improvement tapers past ~80 steps. This pattern suggests that the majority of achievable tasks fall within a moderate difficulty band and that each model approaches a practical ceiling beyond which additional steps yield diminishing returns. The two frontier models, Claude Opus 4.6 and GPT-5.4, reach substantially higher asymptotic perfect score rates (~53% and ~48%, respectively) and climb more steeply through the mid-range, indicating greater per-step efficiency in addition to a higher capability ceiling. Claude Sonnet 4.6 follows a similar trajectory but plateaus lower at ~44%, while GPT-5.4 Mini (~27%) and Qwen 3.5 (35B-A3B and 9B) (~13%) exhibit markedly shallower slopes and lower asymptotes, suggesting that their shortfall reflects fundamental capability gaps rather than insufficient step budgets. In particular, Qwen-3.5, which is one of the best open weights computer-use agents, is significantly worse than frontier API-based models. In future work, it will be valuable to study how we can further improve the long-horizon abilities of these agents, e.g. through reinforcement learning or inference-time search.

Difficulty levels As described in Sec. 3.2, we also categorized Odysseys tasks by their difficulty levels, which are roughly correlated with how many steps are required to solve the task, and how many unique website domains it needs to traverse. We summarize the performance of the models on these tasks in Fig. 4. We observe that performance substantially degrades on hard tasks, with Opus 4.6 averaging a rubric score average of 57.7,

substantially outperforming GPT-5.4 (which achieves a score of 43.0). Opus 4.6 and Sonnet 4.6 models also tend to use more steps across all difficulty levels, particularly on hard tasks.

Failure modes All frontier models we tested exhibit distinct failure patterns on long-horizon tasks. Opus 4.6 most commonly fails through over-investment in research: it continues gathering information but does not transition to producing the required deliverable, such as creating the final document, and eventually exhausts the step budget with an empty output artifact. This pattern appears in 6 of its 12 zero-score runs. This results in Opus runs remaining productive but incomplete until termination, hitting the 100-step cap on 39% of tasks, while its partial-credit runs average 97.4 steps.

GPT-5.4, by contrast, more often fails through inaction despite correct high-level reasoning. On 4 of its 7 zero-score tasks, it generates long and detailed plans that correctly identify which pages to visit and what information to collect, but then terminates with an empty action after only a few steps, sometimes without any browser interaction at all. The two models also share a common failure mode on especially broad tasks with many parallel subtasks. Tasks that require visiting 10 to 30 venues, such as planning a trip to every MLB stadium, often cause both models to become stuck in the first phase of the task, such as collecting schedules, without progressing to later stages like flight search, hotel selection, or document compilation. All such tasks receive zero scores from both models. This suggests that current agents struggle not only with long sequential horizons, but also with high-fanout task structure, where they must allocate effort across many related subtasks. This is one direction for which models would likely benefit from subagents or a multi-agent setup.

Surprising capabilities Both GPT-5.4 and Opus 4.6 models exhibit behaviors that were notably sophisticated. GPT-5.4 in particular develops several strategies that repurpose browser and system primitives for information extraction. We observed it using an unprompted strategy for bulk spreadsheet entry: rather than pasting cell contents directly, it encodes full tables as base64 strings inside its Python action, decodes them at runtime, and pastes the result through the clipboard so that tab and newline delimiters are preserved (Fig. 2). We observe this in 8 runs. Another recurring pattern is that when a product page fails to render correctly due to JavaScript errors, the model navigates to the raw HTML through the view-source: protocol, searches within the source using `ctrl+f`, and extracts structured product metadata from embedded JSON-LD markup, including variant identifiers, stock status, and size mappings (Fig. 8). We observe this behavior in 23 runs.

Opus 4.6 exhibits a different set of strengths. When target websites return 403 errors, it sometimes falls back to the Wayback Machine by constructing archived URLs and retrieving cached versions of otherwise inaccessible pages (Fig. 10). This allows it, for example, to recover class schedule information that would not be available through direct browsing alone. Opus also makes substantially heavier use of middle-click to open links in background tabs without losing its place on the current page, using this interaction 131 times across all runs compared to just 7 for GPT-5.4. In addition, it uses `ctrl+f` not only to locate information, but also to falsify hypotheses about page relevance: after searching for a keyword and observing no matches, it often treats that absence as evidence that the page is unproductive and moves on immediately (Fig. 9).

5 Conclusion

In this work, we introduced Odysseys, a benchmark for evaluating web agents on realistic long-horizon tasks drawn from real browsing behavior. In contrast to previous web benchmarks that emphasize short, single-site episodes, Odysseys focuses on long multi-site workflows that require sustained planning, navigation, and cross-page reasoning. We demonstrated that existing trajectory-level evaluation methods were insufficient for long-horizon trajectories, and proposed rubric-based evaluations which decompose long-horizon success into verifiable intermediate outcomes. This yields substantially higher agreement with human judgment. Our experiments show that, despite strong recent progress in computer-use agents, long-horizon web interaction is far from solved. Even the best frontier models achieve only about 53% perfect-task success, with performance dropping sharply

on harder tasks and plateauing as step budgets increase. Our analysis suggests that current limitations reflect deeper challenges in long-context planning, maintaining coherence across sites, and reliably executing extended workflows. Improving computer-use agents on long-horizon tasks, such as through reinforcement learning or inference-time search, is a promising direction for future work. We release Odysseys tasks, rubrics, and analysis publicly at [removed_for_review](#). We believe the combination of realistic task design and fine-grained evaluation in Odysseys will help drive future progress on computer-use agents that are able to operate robustly over extended time horizons.

Ethics Statement

Intended uses Odysseys is a research benchmark designed to measure and evaluate the progress of computer-use agents on long-horizon web tasks. The models and scaffolds evaluated in this paper are research prototypes assessed in controlled environments, and our results are not intended to endorse their deployment in practical applications at present. Although tasks in Odysseys are executed on the live Internet for the best measure of real world behavior of computer-use agents, the tasks are designed around read-only or low-impact interactions such as browsing, searching, and comparing information, rather than actions with real-world side effects such as completing purchases or submitting forms.

Potential for misuse As computer-use agents become more capable, they could be leveraged for malicious purposes such as automated scraping of personal information, more sophisticated phishing, or circumventing website access controls. We note several emergent behaviors in our analysis in Sec. 4.3, such as using `view-source:` to extract structured metadata from pages that fail to render, or retrieving cached content via the Wayback Machine to bypass access restrictions. These illustrate that agents can repurpose browser primitives in unintended ways. Developers deploying such agents should consider these capabilities and implement appropriate safeguards, including respecting website terms of service and robots.txt policies.

Broader impacts Advancing autonomous web agents has the potential to improve accessibility for users with disabilities or limited technical skills, and to automate tedious, repetitive computer workflows. However, as our results show that even frontier models achieve only 53% perfect task success on Odysseys, and exhibit several key failure modes such as those described above. As their capabilities improve, researchers and developers should carefully consider the economic and social implications, including potential effects on employment. We believe that benchmarks like Odysseys, which expose concrete capability gaps in frontier models, contribute to the responsible development of this technology by enabling the community to measure progress transparently.

Data collection and privacy The browsing history used to construct Odysseys was collected from consenting participants recruited through Prolific, who voluntarily annotated their own Chrome browsing histories. All task prompts were reviewed, and any content that revealed personally identifiable information (PII) was rewritten or removed. No raw browsing histories are released; only the final, de-identified task prompts and rubrics are published.

References

- Peter Anderson, Angel Chang, Devendra Singh Chaplot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, et al. On evaluation of embodied navigation agents. *arXiv preprint arXiv:1807.06757*, 2018.
- Axel Backlund and Lukas Petersson. Vending-bench: A benchmark for long-term coherence of autonomous agents. *arXiv:preprint arXiv:2502.15840*, 2025.
- Rogério Bonatti, Dan Zhao, Francesco Bonacci, Dillon Dupont, Sara Abdali, Yinheng Li, Yadong Lu, Justin Wagle, Kazuhito Koishida, Arthur Bucker, et al. Windows agent arena: Evaluating multi-modal os agents at scale. *ICML*, 2025.
- Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. Seeclick: Harnessing gui grounding for advanced visual gui agents. *ACL*, 2024.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web. *NeurIPS*, 2023.
- Alexandre Drouin, Maxime Gasse, Massimo Caccia, Issam H Laradji, Manuel Del Verme, Tom Marty, Léo Boisvert, Megh Thakkar, Quentin Cappart, David Vazquez, et al.

- Workarena: How capable are web agents at solving common knowledge work tasks? *ICML*, 2024.
- Google. Finding answers gets better with chrome. <https://blog.google/products-and-platforms/products/chrome/finding-answers-gets-better-chrome/>, 2022.
- Anisha Gunjal, Anthony Wang, Elaine Lau, Vaskar Nath, Yunzhong He, Bing Liu, and Sean Hendryx. Rubrics as rewards: Reinforcement learning beyond verifiable domains. *ICLR*, 2026.
- Helia Hashemi, Jason Eisner, Corby Rosset, Benjamin Van Durme, and Chris Kedzie. Llm-rubric: A multidimensional, calibrated approach to automated evaluation of natural language texts. *ACL*, 2024.
- 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. *ACL*, 2024.
- Raghav Kapoor, Yash Parag Butala, Melisa Russak, Jing Yu Koh, Kiran Kamble, Waseem Alshikh, and Ruslan Salakhutdinov. Omniact: A dataset and benchmark for enabling multimodal generalist autonomous agents for desktop and web. *ECCV*, 2024.
- Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Russ Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks. *ACL*, 2024.
- Thomas Kwa, Ben West, Joel Becker, Amy Deng, Katharyn Garcia, Max Hasin, Sami Jawhar, Megan Kinniment, Nate Rush, Sydney Von Arx, Ryan Bloom, Thomas Broadley, Haoxing Du, Brian Goodrich, Nikola Jurkovic, Luke Harold Miles, Seraphina Nix, Tao Lin, Neev Parikh, David Rein, Lucas Jun Koba Sato, Hjalmar Wijk, Daniel M. Ziegler, Elizabeth Barnes, and Lawrence Chan. Measuring ai ability to complete long software tasks. *arXiv:preprint arXiv:2503.14499*, 2026.
- Kaixin Li, Ziyang Meng, Hongzhan Lin, Ziyang Luo, Yuchen Tian, Jing Ma, Zhiyong Huang, and Tat-Seng Chua. Screenspot-pro: Gui grounding for professional high-resolution computer use. *MM*, 2025.
- Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. Reinforcement learning on web interfaces using workflow-guided exploration. *ICLR*, 2018.
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. *ICLR*, 2024.
- Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia: a benchmark for general ai assistants. *ICLR*, 2024.
- OpenAI. text-embedding-3-small. <https://developers.openai.com/api/docs/models/text-embedding-3-small>, 2024.
- Christopher Rawles, Sarah Clinckemauille, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Marybeth Fair, Alice Li, William Bishop, Wei Li, Folawiyi Campbell-Ajala, et al. Android-world: A dynamic benchmarking environment for autonomous agents. *ICLR*, 2025.
- Rulin Shao, Akari Asai, Shannon Zejiang Shen, Hamish Ivison, Varsha Kishore, Jingming Zhuo, Xinran Zhao, Molly Park, Samuel G Finlayson, David Sontag, et al. Dr tulu: Reinforcement learning with evolving rubrics for deep research. *arXiv preprint arXiv:2511.19399*, 2025.
- Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Percy Liang. World of bits: An open-domain platform for web-based agents. *ICML*, 2017.

- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Musique: Multihop questions via single-hop question composition. *TACL*, 2022.
- Vijay Viswanathan, Yanchao Sun, Shuang Ma, Xiang Kong, Meng Cao, Graham Neubig, and Tongshuang Wu. Checklists are better than reward models for aligning language models. *NeurIPS*, 2025.
- Weixuan Wang, Dongge Han, Daniel Madrigal Diaz, Jin Xu, Victor Rühle, and Saravan Rajmohan. Odysseybench: Evaluating llm agents on long-horizon complex office application workflows. *arXiv:preprint arXiv:2508.09124*, 2025.
- Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. Browsecomp: A simple yet challenging benchmark for browsing agents. *arXiv:preprint arXiv:2504.12516*, 2025.
- Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. Travelplanner: A benchmark for real-world planning with language agents. *ICML*, 2024a.
- Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh J Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *NeurIPS*, 2024b.
- Tianbao Xie, Jiaqi Deng, Xiaochuan Li, Junlin Yang, Haoyuan Wu, Jixuan Chen, Wenjing Hu, Xinyuan Wang, Yuhui Xu, Zekun Wang, Yiheng Xu, Junli Wang, Doyen Sahoo, Tao Yu, and Caiming Xiong. Scaling computer-use grounding via user interface decomposition and synthesis. *NeurIPS*, 2025.
- Frank F. Xu, Yufan Song, Boxuan Li, Yuxuan Tang, Kritanjali Jain, Mengxue Bao, Zora Z. Wang, Xuhui Zhou, Zhitong Guo, Murong Cao, Mingyang Yang, Hao Yang Lu, Amaad Martin, Zhe Su, Leander Maben, Raj Mehta, Wayne Chi, Lawrence Jang, Yiqing Xie, Shuyan Zhou, and Graham Neubig. Theagentcompany: Benchmarking llm agents on consequential real world tasks. *ICML*, 2025.
- Tianci Xue, Weijian Qi, Tianneng Shi, Chan Hee Song, Boyu Gou, Dawn Song, Huan Sun, and Yu Su. An illusion of progress? assessing the current state of web agents. *COLM*, 2025.
- Pei Yang, Hai Ci, and Mike Zheng Shou. Macosworld: A multilingual interactive benchmark for gui agents. *arXiv preprint arXiv:2506.04135*, 2025.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *EMNLP*, 2018.
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. *NeurIPS*, 2022.
- Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. τ -bench: A benchmark for tool-agent-user interaction in real-world domains. *ICLR*, 2025.
- Yutori. Introducing navigator. <https://yutori.com/blog/introducing-navigator>, 2025. Yutori Blog.
- Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for building autonomous agents. *NeurIPS*, 2024.

A Appendix

A.1 Large Language Model Disclosure

We iteratively use a coding agent (Claude Code) with a human in the loop to generate several plots in this paper conditioned on numerical results and visualization instructions. We also used Claude Code to flag potentially interesting examples in the trajectories and write a description of these examples. We reviewed these outputs and polished these for our qualitative analysis section of the paper. We also give Claude Code access to our codebase, results, and data, to flag any issues with our paper and provide feedback for iterative writing.

A.2 Data Collection Interface

Figure 5 shows the desktop application used by participants to annotate their Chrome browsing journeys (Sec. 3.2). For each journey, the interface displays the segmented URLs and guides the participant through four annotation steps: (1) selecting the key URL that represents the task’s success state, (2) indicating an automation preference, (3) writing a task label as they would prompt an AI tool, and (4) judging feasibility.

Journey steps

To upload a journey, the following four items must be completed:

- 1) Pick the single URL from the journey steps that represents your success state and click "Save key URL"
- 2) Would you want this activity to be automated or done for you?
- 3) Task Question
- 4) Feasibility

pittsburgh new restaurants - Google Search
2026-03-28T22:11:02.343541 - https://www.google.com/search?q=pittsburgh+new+restaurants&rlz=1C5CHFA_enUS1106US1106&ooq=pittsburgh+new+res&gs_lcrp=EgZjaHJvbWUqCQgBEAAyDRIABDIGCAAQRrg5MgkIARAAGA0YgAQyCQgCEAA8

pittsburgh new restaurants - Google Search
2026-03-28T22:11:11.115650 - https://www.google.com/search?q=pittsburgh+new+restaurants&rlz=1C5CHFA_enUS1106US1106&ooq=pittsburgh+new+res&gs_lcrp=EgZjaHJvbWUqCQgBEAAyDRIABDIGCAAQRrg5MgkIARAAGA0YgAQyCQgCEAA8&liq=ChpwaxR0c2J1cmdoIG5ldyByZXN0YXVyYW50c0UlpuNjI2AgAhalBAGAAIFnBpdHRzYnVyZ2ggcmVzdGF1cmFudHMqAggDkgEKcmVzdGF1cmFudOABAA#rlimm=21859563

Spring Into These New (and Soon-to-Open) Restaurants, Cafes and Breweries | Pittsburgh Magazine
2026-03-28T22:11:17.877250 - https://www.pittsburghmagazine.com/spring-2026-restaurant-openings/

Step 1: Pick the one URL that best represents the success state for this journey, then click "Save key URL" (required before uploading).

Pick the URL that shows you achieved your goal, then click "Save key URL" to lock it in.

2) Would you want this activity to be automated or done for you?

Yes - I'd want this activity automated or done for me

3) Task Question

Write a request to someone that did the following activity. Be as descriptive as possible, as you would prompt an AI tool.

Minimum 15 characters.

Good example: Find the cheapest used Honda Civic that meets all the following criteria: under \$25,000, fewer than 30,000 miles, includes Apple CarPlay, adaptive cruise control, and blind-spot sensors.

Bad example: Find car.

4) Feasibility

Were you able to find the information or do what you wanted to?

Feasible

Figure 5: The annotation interface used by participants to label their Chrome browsing journeys. The interface guides participants through four steps: selecting the key URL representing the success state, indicating an automation preference, writing a descriptive task label, and judging feasibility.

A.3 Journey Refinement

The 2,380 raw journeys collected from participants required substantial refinement before they could serve as building blocks for Odysseys tasks.

LLM screening An LLM screens each journey along four dimensions: (1) label accuracy (whether the task description matches what the visited URLs actually support), (2) feasibility (whether the task can be completed without login credentials or additional personal context), (3) login requirement (whether any step requires authentication), and (4) overall quality. For each journey, the LLM also generates a refined label scoped to what the URLs support, along with notes justifying each assessment. This process judges 71.8% of the original labels as inaccurate relative to the URLs. This high rate reflects a combination of factors: participants sometimes described their *intent* rather than what the URLs show (e.g., “find a good laptop” when the URLs only show a single product page), submitted labels for browsing sessions that spanned multiple unrelated goals, or wrote underspecified descriptions that could not be verified against the visited pages. The LLM-refined labels address these issues by anchoring the task description to the observable URL evidence.

Manual review Each journey is then manually reviewed by the authors through the annotation interface shown in Appendix A.2. Reviewers independently verify URL selections, assess feasibility by manually searching the web, and choose between the original participant label, the LLM-refined label, or a custom label written by the reviewer. After filtering for feasibility, removing tasks that require login-gated flows, and discarding low quality tasks, 696 usable journeys (29.2% of the original 2,380) remained.

A.4 Task Composition Pipeline

The 696 refined journeys are short, single-site subtasks. To compose them into long-horizon Odysseys tasks, we first cluster related journeys, then use an LLM to chain subsets into coherent multi-step workflows.

Clustering We embed each journey label using text-embedding-3-small (OpenAI, 2024), reduce the embeddings to 15 dimensions with UMAP, and cluster them using HDBSCAN. UMAP preserves both local and global embedding structure, and HDBSCAN automatically determines the number of clusters, accommodates variable-density groups, and assigns outlier journeys to a noise class rather than forcing them into ill-fitting clusters. We then construct a *theme graph* by connecting clusters whose centroid cosine similarity lies within a fixed range, filtering out both weakly related and near-duplicate clusters.

Chaining Multi-step Odysseys tasks are composed by traversing this theme graph. For each task, we select a seed cluster at random, run a 1–2 hop BFS to collect 3–6 related clusters, and draw 3–4 representative journeys per cluster to form 12–24 candidate journeys. GPT-5.4 then selects and orders a subset into a coherent workflow, generating a natural-language prompt, step plan with transitions, rubric with verification procedures, and a coherence score in a single call.

Composition prompt The system prompt below defines the task format and requirements (few-shot style examples are omitted for brevity):

```
You are an expert at designing realistic, long-horizon web agent benchmark tasks.

You will receive a MENU of atomic web browsing tasks. Each has an ID, website,
description, cluster theme, and optional location. Your job is to compose a
SUBSET into a single coherent multi-step task that a real person would actually do.

CRITICAL REQUIREMENTS:
1. SEQUENTIAL FLOW: Each step must logically lead to the next. Later steps must
   USE or BUILD ON information from earlier steps.
2. UNIFIED GOAL: All steps serve one overarching purpose.
3. GEOGRAPHIC CONSISTENCY: If steps involve physical locations, they must be in
   the same city/region.
4. CROSS-SITE: Use at least 2 different websites.
5. INFORMATION DEPENDENCIES: At least 30% of steps should depend on a prior
   step's output.
6. NATURAL VOICE: The agent_prompt must sound like a real person talking to an
   assistant -- conversational, with personal context.

STYLE EXAMPLES: [8 few-shot examples spanning easy -> very_hard]

OUTPUT FORMAT: { goal, task_name, agent_prompt, steps, rubric, dependencies,
                 skills, primary_skill, deliverable, self_score, reasoning }
- Select exactly {target_length} steps for the task.
- You may SKIP menu items that don't fit.
- If you cannot form a coherent task, return {"steps": [], "self_score": 0}.
```

The user prompt provides the target difficulty, step count, and a menu of candidate journeys drawn from the related clusters:

```
Target: 6 steps, hard difficulty

Available tasks:
[A0] [google.com] Search for round-trip flights from Pittsburgh to LAX
      cluster: flight search and booking
[A1] [booking.com] Find hotel near airport with overnight stay option
      cluster: hotel comparison and booking
[A2] [maps.google.com] Check drive time between two locations
      cluster: route planning and navigation
[A3] [enterprise.com] Search for car rental options at airport
      cluster: car rental comparison
...

Compose a coherent 6-step sequential workflow.
```

A.5 Odysseys Task QA Interface

After composition and prompt rewriting (Sec. 3.3), every chained Odysseys task is reviewed by the authors using the QA interface shown in Fig. 6. The interface lists all tasks with their difficulty level, rubric count, and score. Expanding a task reveals the full prompt and individual rubric items with weights, requirements, and verification criteria, allowing reviewers to verify coherence and actionability before the task enters the benchmark.

The screenshot displays the Odysseys QA interface. At the top, task #0 is expanded, showing its title, difficulty level (easy), rubric count (4), and CUA score (7). Below the title is a text area containing the prompt: "I'm putting together a small TV watchlist and want to anchor it around The Pitt first, so please go to Hulu and open the actual show page for The Pitt to confirm what service it's on, then leave that tab open so I can see the listing myself. Once you've confirmed that, use Wikipedia to look up the TV series Ponies and pull the main cast names from the series page so I can compare who's in a different show; if the cast is listed on the page, open the Ponies article itself and keep that tab available too just so I can glance at it. Then round out the watchlist with something older by going to Memory Alpha and finding the entry for Amok Time, and grab the key details from that page including which Star Trek series it belongs to, the season and episode number, and the original air date. Please give me everything back in one concise summary with the streaming service for The Pitt, the Ponies cast list, and the Amok Time details, and keep the Hulu and Memory Alpha pages open in separate tabs so I have visual proof." Below the prompt is a table of rubrics with 4 items, each with an ID, weight, requirement, and verification criteria. A "Save this task" button is visible below the rubric table. Below task #0, four other tasks are listed in a collapsed state, each showing its title, difficulty level, rubric count, and CUA score.

ID	WEIGHT	REQUIREMENT	VERIFICATION
R1	0.3	The agent opens the Hulu page for "The Pitt," confirms the serv	A grader can see the Hulu title page for "The Pitt" open in a bro
R2	0.3	The agent extracts the main cast members for the TV series "Pc	A grader can inspect the open Wikipedia article for "Ponies" and
R3	0.25	The agent reports the Memory Alpha details for "Amok Time," in	A grader can view the open Memory Alpha page for "Amok Time
R4	0.15	The agent provides one concise combined summary covering th	The final response is a single concise summary containing all th

Figure 6: The Odysseys QA interface used for manual review of chained tasks. Reviewers can expand each task to inspect the full prompt and individual rubric items with weights, requirements, and verification criteria.

A.6 Full Task Descriptions for Table 2

The following are the complete, unabridged task prompts for the six representative tasks shown in Table 2, listed in order of difficulty. These are the exact instructions given to the computer-use agent.

Easy: Kitchen setup (Brussels sprouts, Le Creuset, dishwasher air gap)

I'm setting up a new kitchen and want one of the first things I make to be roasted Brussels sprouts, so could you start on Google and find me a recipe that clearly uses both Parmesan and balsamic vinegar, then open the actual recipe page and note the title, oven temperature, and cook time because I want to make sure the cookware I buy fits that kind of roasting setup. Once you've got that recipe open, head to Le Creuset and look for a light green Dutch oven, and specifically check whether the 5.5 qt size is offered in that color so I can see if it would work for recipes in that range; please open the product page itself and leave it open so I can look at the color and size options on the page. While you're at it, I'm also sorting out kitchen appliances before I start cooking, so go to YouTube and find a practical video about whether a dishwasher installation needs an air gap, open the video page, start playing it, and tell me the main decision points like when an air gap is required, when a high loop is used instead, and why the air gap exists in the first place. Please keep the recipe tab, the Le Creuset product tab, and the YouTube video tab open in separate tabs so I can compare everything visually afterward.

Easy: Weather risk snapshot (Baltimore, Syracuse, Rittman, Mount Holly)

I'm trying to decide whether driving this week is a bad idea, so can you build me a quick weather risk snapshot that starts with what it feels like right now and then widens out to the bigger trouble spots? First, on Google, search for Baltimore, Maryland weather and grab the current temperature plus the plain-English condition like cloudy, sunny, rain, or whatever it says, just so I have a baseline for home conditions. Then go to Wunderground and look up Syracuse, New York, and check the 10-day/7-day style forecast to find the lowest temperature expected over the next 7 days, including which day it happens, because I want to compare that colder destination against Baltimore. After that, use the National Weather Service forecast page for the Rittman, Ohio area near Marshallville and tell me what the current forecast says and whether there are any active alerts posted there, since that would really affect an Ohio leg of the drive; please open the actual forecast page and leave it visible so I can see the alert area and forecast text myself. Finally, go to the NWS Mount Holly page, find the winter forecast graphic, and report the snowfall amount shown there so I can tell whether the Mid-Atlantic part looks like a nuisance event or something more serious; if the graphic opens separately, leave that tab open too so I can look at the map. In the end, send me a short location-by-location summary with the key weather risk for Baltimore, Syracuse, Rittman/Marshallville, and the Mount Holly region.

Medium: Japan trip planning for a Canadian in the U.S.

I'm a Canadian citizen living in Pittsburgh, PA, and my passport expires in about 3 months, so I'm trying to get everything sorted before a 2-week tourist trip to Japan. Could you start on the official Government of Canada site and find the passport renewal process for a Canadian living in the U.S., including the exact renewal form I'd need, the supporting documents, photo rules, whether I need a guarantor or references, the fee in CAD, how I'm supposed to submit it from the U.S., and the current processing time, because I need to know if this is realistic before I book anything. Once you have that, use the official Canadian embassy/consulate pages to figure out which Canadian mission is closest to Pittsburgh, Pennsylvania 15222 that handles passport services, and open the actual office page so I can see the address, passport service hours, and whether I need an appointment or have to use some booking request process; please leave that page open. After that, check Japan's official Ministry of Foreign Affairs site to confirm whether a Canadian passport holder going to Japan for tourism for 2 weeks needs a visa, and note any conditions or exceptions that matter. Then go to the Government of Canada travel advisory page for Japan and tell me the current advisory level plus any

highlighted health, safety, or entry notes, and keep that advisory page open in another tab so I can look at it myself. Finally, compare travel insurance options on PolicyAdvisor.com and Kanetix.ca for this situation: a Canadian citizen currently living in the U.S. who wants coverage connected to travel to Japan, and I mainly want to see whether either site shows plans that would actually work for someone based in the U.S. rather than Canada, so please capture provider names, medical emergency coverage, trip cancellation/interruption if shown, and any residency or eligibility restrictions. If either site has useful quote or results pages, open the most relevant options in separate tabs so I can compare them visually. At the end, give me a concise summary that ties all of this together and clearly points out any uncertainty, especially around insurance eligibility for a Canadian living in the U.S.

Medium: Tesla Model 3 lease evaluation in Los Angeles

I'm trying to sanity-check whether moving ahead with a Tesla Model 3 lease in Los Angeles is actually manageable month to month, so start on Google and look up current Tesla Model 3 lease pricing for the Los Angeles area, including the lease term, due-at-signing amount, and any discounts, tax credits, or rebates you can find, because I want a realistic baseline instead of just a headline number. Once you've got that monthly lease figure, use it as a reference point and go to the Los Angeles Craigslist site, specifically the San Gabriel Valley section, and find at least three trailer listings that look like plausible live-in fallback options under \$10,000; open the actual posting pages in separate tabs so I can see the photos and verify the listings are still live, then note each one's title, price, and location. After that, go to Zillow and look for an LA-area rental whose monthly price is in the same ballpark as the Tesla lease payment you found, so I can compare whether paying for housing at that level makes more sense than taking on the car; open the actual Zillow listing page and leave it open so I can check the photos and map myself. In the end, give me a short comparison that includes the Tesla lease deal, the three Craigslist trailer backups, and the Zillow rental option that's closest in monthly price to the lease.

Hard: 30 MLB stadiums summer road trip

I'm daydreaming about doing a ridiculous-but-fun summer baseball trip where I see exactly one game at all 30 MLB stadiums, and I want you to build the whole thing in a way I could actually use. Start on MLB.com and pull the official summer schedule so we can choose one real game date at each stadium, and please lean toward matchups where I might get to see stars I care about most like Shohei Ohtani, Aaron Judge, and Ronald Acuna Jr. whenever that's realistically possible. As you're picking games, open the actual game or team schedule pages in separate tabs for a few representative stops so there's visible proof the dates are live, and keep the key schedule tabs open so I can glance at them later. Once you've got the 30 stadium/date choices, use Google Flights to figure out the smartest sequence between stops and compare flights versus driving for each leg, using whatever is cheaper and more practical in summer, because I want this to feel like a real budget-conscious trip instead of fantasy routing. After that, use Booking.com to find one hotel option for each game night that's reasonably close to the stadium—something like within about 2 miles if possible and not outrageously priced for a solo traveler—and open at least a couple of the actual hotel listing pages with photos/maps so I can visually sanity-check the neighborhoods. Then use Yelp to find at least one must-try local food spot near each stadium, ideally something iconic to that city or ballpark area, and open a few of the restaurant pages so I can see that they're real places with reviews. Finally, put everything into a CryptPad Document in one organized itinerary with each stadium listed exactly once, the chosen game and matchup, whether it includes Ohtani, Judge, Acuna, or another notable player, the travel leg before it with the cheaper mode and estimated cost, one hotel with estimated nightly price, one food pick, and running totals so I can see what this insane summer would actually cost. Leave the finished CryptPad Document open at the end, and if you create any comparison tabs along the way, keep the most useful ones open so I can review them.

Hard: Banana bread recipe synthesis

I want to develop the best banana bread recipe. Look up the top 10 recipes online (by engagement, popularity, reviews) and compare the recipes (e.g. composition of ingredients, additions, cooking method), identifying and highlighting similarities and unique points that make each recipe good. Keep the most unique or highly reviewed 3 recipes in open tabs so I can reference them, and make sure at least one has a YouTube video (also keep this video open and start playing it). Then, from these three, create the best recipe you can combining aspects of these and provide me with step by step instructions.

A.7 Human Agreement Annotation Interface

Figure 7 shows the trajectory viewer used by annotators to evaluate rubric satisfaction for each agent run (Sec. 4.1). The interface displays the agent’s step-by-step actions alongside screenshots, reasoning, and task rubrics. For each rubric item, the annotator marks whether it was satisfied during the trajectory. Ambiguous cases are flagged for discussion.

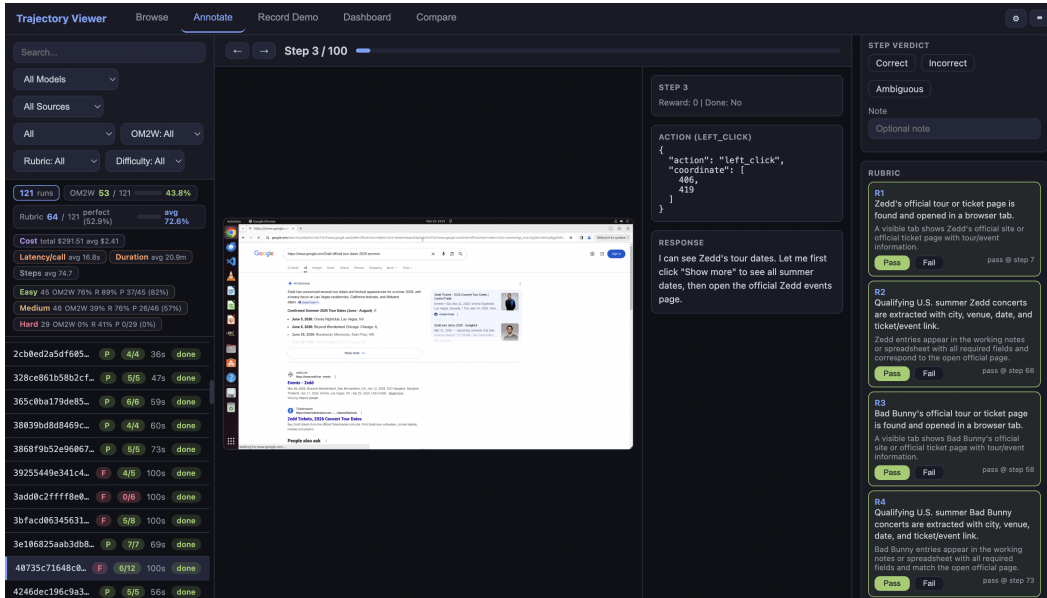


Figure 7: The trajectory viewer used for human agreement annotation. Annotators step through the agent’s actions (left panel), view screenshots and action details (center), and judge each rubric item as pass or fail (right panel).

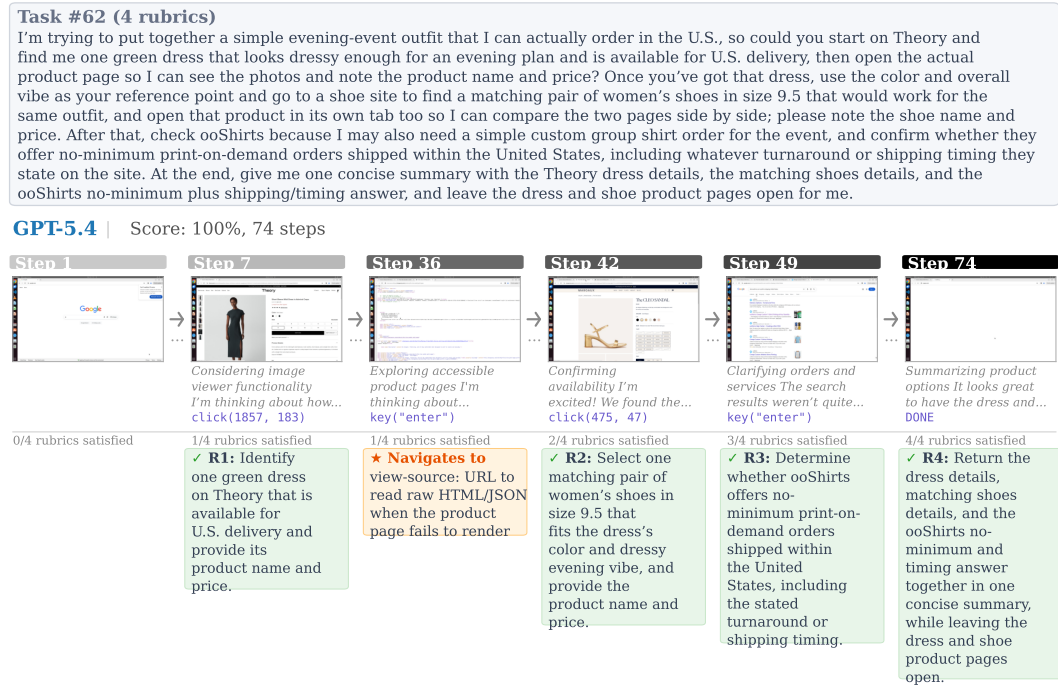


Figure 8: When margauxny.com rendered blank due to JavaScript failures, the GPT-5.4 agent typed view-source:https://margauxny.com/products/... directly in the address bar, then used ctrl+f to search for variants, InStock, 40.5, and US 9.5 in the raw HTML. From the embedded JSON-LD schema, it decoded EU 40.0 = US 9.5, confirmed that the SKU was in stock, and reconstructed the direct variant URL, all without the product page ever visually rendering.

A.8 Further Qualitative Results

The following figures illustrate the surprising agent behaviors discussed in Sec. 4.3.



Figure 9: After navigating to `brave.com/linux`, Opus used `ctrl+f` to search for Chromebook and observed `0/0` matches, confirming that the page did not cover the topic at all, and immediately pivoted to a different strategy. Rather than scrolling to verify, it treated the absence of a match as decisive evidence that the page was unhelpful.

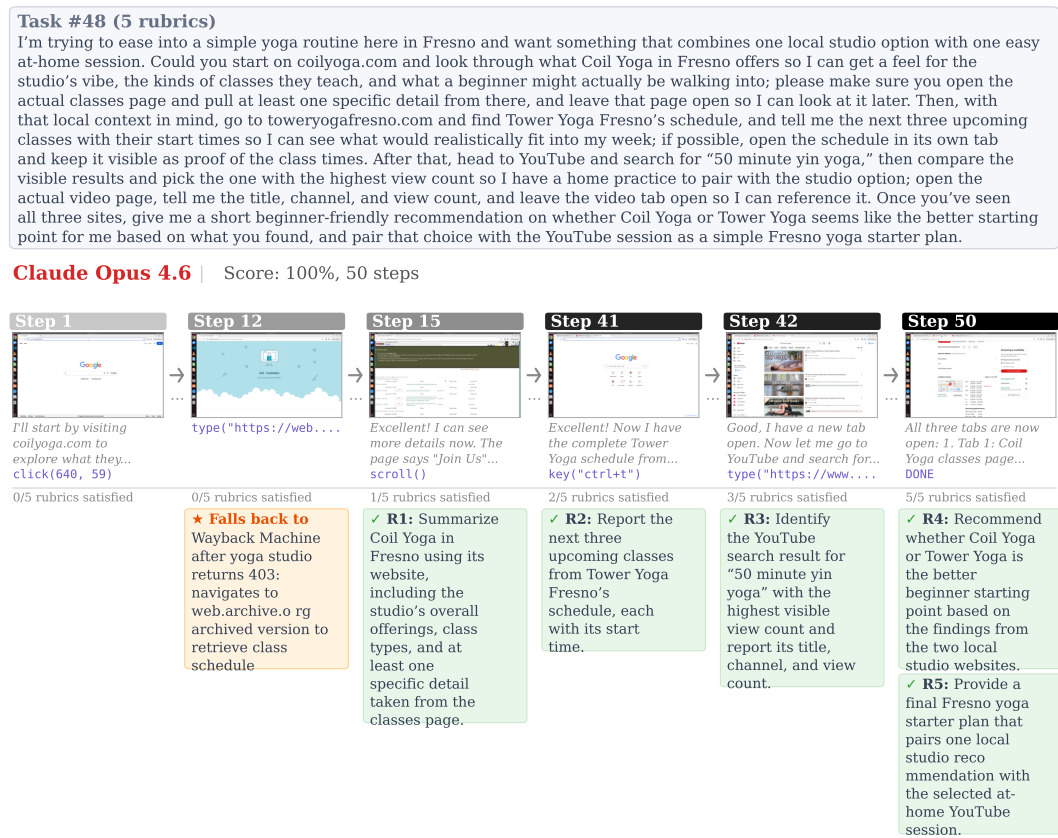


Figure 10: Both yoga studio websites returned 403 Forbidden. Opus immediately navigated to <https://web.archive.org/web/2024/https://coilyoga.com/classes/> and repeated the same strategy for Tower Yoga, successfully retrieving archived pages from June 2024 that contained the full class schedule information.