

TOWARDS INTERNET-SCALE TRAINING FOR AGENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

The predominant approach for training web navigation agents is to gather human demonstrations for a set of popular websites and hand-written tasks, but it is becoming clear that human data is an inefficient resource. We develop a pipeline to facilitate internet-scale training for agents without laborious human annotations. In the first stage, an LLM annotates 150k sites with agentic tasks. In the next stage, LLM agents complete tasks and produce trajectories. In the final stage, an LLM filters trajectories by judging their success. Language models are powerful data curation tools, identifying harmful content with an accuracy of 97%, judging successful trajectories with an accuracy of 82.6%, and producing effective data. We train agents based on *Qwen 3 1.7B* that are competitive with frontier LLMs as web agents, while being smaller and faster. Our top agent reaches a success rate of 56.9%, outperforming the data collection policy *Qwen 3 235B*, a 235 times larger *Llama 4 Maverick*, and reaching 94.7% of the performance of *Gemini 2.5 Flash*. We will be releasing code, models and data to reproduce the entire pipeline.

1 INTRODUCTION

The predominant approach for training LLM web navigation agents is to collect human demonstrations for a set of manually curated websites and tasks (Deng et al., 2023; Zhou et al., 2024b; Putta et al., 2024; Koh et al., 2024a; Liu et al., 2024; Lù et al., 2024; Rawles et al., 2023). Human data can be laborious to collect, and becomes costly to scale as the breadth of skills that users require from language model agents grows. There are more than 300M sites on the western internet according to The Common Crawl Foundation (2025), and the range of sites that researchers have annotated represents a tiny fraction of the vast available data. And crucially, the existing human data is *static*. There is a growing need to automate pipelines for training the next generation of language model agents in a *dynamic internet-scale environment*. This paper addresses the core challenge of building this environment—reducing dependence on human annotations. We develop an automatic pipeline that aims to facilitate internet-scale training for agents, which we refer to as the InSTA pipeline.

The InSTA pipeline has three stages. In the first stage, we employ a language model task proposer that annotates 150k sites with live web navigation tasks for agents to perform. Existing works are limited to 200 popular websites (He et al., 2024; Lù et al., 2024; Rawles et al., 2023; Deng et al., 2023; Zhou et al., 2024c; Murty et al., 2025) that researchers have annotated manually, and a handful of synthetic websites (Zhou et al., 2024b; Koh et al., 2024a; Yao et al., 2023a). Our first goal in this paper is to improve coverage of real-world sites. To accomplish this, we cast a wide net. Starting from the top 1M sites on the internet ranked by popularity, our task proposer filters down to 150k sites that have safe content. Safety is critical when building autonomous agents, and our task proposer succeeds at detecting harmful content with an accuracy of 97%. Tasks are generated for sites marked as safe by the task proposer, and we run language model agents to complete the generated tasks. Agent progress is then fed back to the task proposer, which reviews trajectories and judge evaluations in order to assign a harder task based on the latest content of the website, closing the loop.

Scaling the task generation loop, we annotate 150k diverse sites with challenging agentic tasks, and release an official huggingface dataset: *data-for-agents/insta-150k-v2*. Motivated by the importance of internet data to progress in modern deep learning, our second goal in this paper is an internet-scale data flywheel for training LLM agents. We approach this by harnessing LLMs as data curation tools. After generating tasks, the pipeline employs pretrained LLMs as agents to complete tasks and produce trajectories, which are filtered by a judge to select the best data. Agents control a virtual web browser and produce reasoning traces that contain function calls to interact with and navigate live webpages.

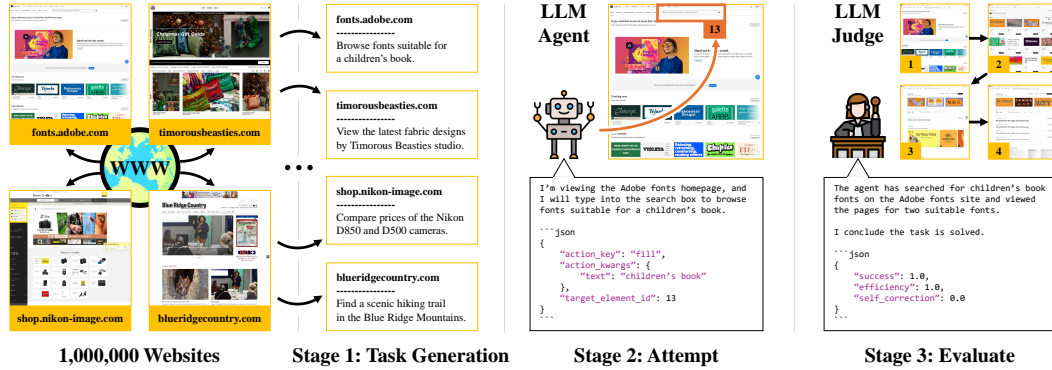


Figure 1: **Overview of the InSTA pipeline.** Our work unlocks a dynamic internet-scale environment that allows training small models to match frontier LLMs as agents, on a fraction of the budget. Starting from the top 1M sites on the internet, we annotate 150k sites with challenging agentic tasks, and release the entire pipeline, including code, models and an official huggingface dataset, on our website: anonymous-insta-pipeline.github.io.

The judge produces a reasoning trace that considers whether a trajectory is successful, and scores the agent on a continuous scale from 0 to 1. We scale the pipeline to create a large reasoning dataset for multimodal agents, with 2.2M screenshots, 2.2M traces for actions, and 150k traces for the judge.

Our data unlocks great potential in small language models as agents. We train a series of models based on *Qwen 3 1.7B* on varying scales of data from the InSTA pipeline, and match the performance of frontier LLM agents, on a fraction of the budget. Our top agent has a success rate of 56.9%, outperforming the data collection policy *Qwen 3 235B*, beating a 235 times larger *Llama 4 Maverick*, and reaching 94.7% of the performance of *Gemini 2.5 Flash*. To share our progress, we are releasing the entire pipeline, including code, models and data, on our website: anonymous-insta-pipeline.github.io.

2 RELATED WORKS

Language Model Agents. There is an emerging paradigm in modern NLP using language models (Radford et al., 2019; Brown et al., 2020; Touvron et al., 2023a;b) as the backbone for agents (Andreas, 2022). These models show impressive reasoning capabilities (Bubeck et al., 2023; Zhong et al., 2024; Valmeekam et al., 2024) that allow them to generalize to downstream applications, such as web navigation, where observations differ from LLM training data. Search algorithms provide a secondary axis to improve the reasoning capabilities of language model agents (Yao et al., 2023b; Besta et al., 2024; Koh et al., 2024b; Zhou et al., 2024a) by providing an explicit algorithmic scaffold and allowing test-time compute to improve reasoning steps (Snell et al., 2024; Zhong et al., 2024). Although most recent work focuses on running language models as zero-shot agents, fine-tuning language models to improve their effectiveness as agents is becoming popular (Putta et al., 2024; Zeng et al., 2023; Zhang et al., 2023; Hong et al., 2023; Xie et al., 2024; Wang et al., 2024; Zhou et al., 2024c; Murty et al., 2025) as target benchmarks are becoming more difficult to solve zero-shot.

Agent Pipelines. There are a growing number of agent pipelines aimed at fine-tuning language models to improve their effectiveness as agents (Mitra et al., 2024; Zeng et al., 2023; Putta et al., 2024; Chen et al., 2023; Ou et al., 2024; Zhou et al., 2024c; Murty et al., 2025). However, driven by the limited data available, many such works train on data with significant overlap with their test environment, either with different tasks for the same environment configuration as the test setting (Deng et al., 2023; Zhou et al., 2024c; Murty et al., 2025), or even the same tasks (Putta et al., 2024). We consider a setting where tasks and environment configurations (i.e. websites) are entirely separate between training and testing, creating a strong train-test split. This presents a challenge: human data for training LLM agents is limited (Deng et al., 2023; Lù et al., 2024). We address this challenge by reducing dependence on human data in agent pipelines. We employ LLMs as data curation tools to automatically design challenging tasks, and select the best training data. Our pipeline allows us to train small models that match frontier LLMs on Web Voyager (He et al., 2024), without using any data from Web Voyager. Contrast this with methods that train primarily on data from Web Voyager (Zhou et al., 2024c; Murty et al., 2025), and may not transfer to other benchmarks (Xue et al., 2025).

Agent Datasets. Datasets for training web navigation agents typically rely on human annotators to create tasks (Zhou et al., 2024b; Koh et al., 2024a; Rawles et al., 2023), and provide demonstrations for tasks (Deng et al., 2023; Lù et al., 2024; Rawles et al., 2023; Shen et al., 2024). However, the amount of data researchers have annotated represents a tiny fraction of the available internet data. There are more than 300M sites on the internet according to The Common Crawl Foundation (2025), yet existing datasets are limited to 200 popular sites that human annotators are familiar with (Deng et al., 2023; Lù et al., 2024; Shen et al., 2024). Human data can be laborious to collect, and becomes costly to scale as the capabilities users require from agents grows. And crucially, human data is *static*. Our work moves away from fixed datasets for training agents, and towards a dynamic internet-scale environment that grows with an ever-changing internet. We are not the first to build an environment (Zhou et al., 2024b; Koh et al., 2024a; Yao et al., 2023a; He et al., 2024), nor are we the first to consider synthetic data (Gandhi et al., 2024; Ou et al., 2024; Setlur et al., 2024; Tajwar et al., 2024), but we have solved a key challenge that unlocks the internet as the largest environment for agents.

Language Model Judges. Using LLMs to judge the correctness of responses is becoming popular to refine LLM predictions (Li et al., 2024), including to verify reasoning steps (Zhang et al., 2024), rejection sample (Snell et al., 2024; Sun et al., 2024), prioritize frontier nodes in search algorithms (Zhou et al., 2024a; Koh et al., 2024b), filter out harmful responses (Inan et al., 2023), provide feedback for response improvement (Madaan et al., 2023; Paul et al., 2024; Patel et al., 2024; Yuksekgonul et al., 2024), and provide ratings for alignment (Lee et al., 2024; Ouyang et al., 2024). Our use of language models to evaluate agents is inspired by Generative Verifiers (Zhang et al., 2024), and the multimodal verifier in He et al. (2024). One difference is our verifier predicts a reasoning trace that scores the agent from 0 to 1, which helps us rank trajectories to select the best data.

3 LANGUAGE MODEL AGENTS

Language model agents are a class of decision-making agents represented by $\pi_{\text{LLM}}(\mathbf{a}_t | \mathbf{s}_t, \mathbf{c}_n)$, a policy that processes multimodal observations \mathbf{s}_t (from a virtual web browser in our case) and predicts textual actions \mathbf{a}_t to complete a task \mathbf{c}_n . Underneath this abstraction, a large language model (LLM) generates actions via the next-token prediction, conditioned on a system prompt $\mathbf{x}_{\text{agent}}$.

$$\mathbf{a}_t = f^{\text{text} \rightarrow \text{act}}(\text{LLM}(\mathbf{x}_{\text{agent}}, \mathbf{c}_n, \text{Enc}(\mathbf{s}_t))) \quad (1)$$

Environment representations for observations and actions typically differ from the expected input format of the language model (typically images and text), and functions are introduced that map the observations to a multimodal prompt $\text{Enc}(\cdot)$, and parse actions from the language model generated response $f^{\text{text} \rightarrow \text{act}}(\cdot)$. For web navigation, the environment state \mathbf{s}_t is HTML DOM, and is often formatted as raw HTML code, an Accessibility Tree, Set-of-marks, or screenshots (Zhou et al., 2024b; Koh et al., 2024a; Chezelles et al., 2024; Shen et al., 2024). We built a fast Markdown parser that converts webpage observations into a compact readable format (refer to the code). Action formats vary between works, and we build on Schick et al. (2023)’s function-calling framework, where a language model generates code that is parsed into a function name and corresponding arguments. Given sets of function and argument names L_i , and sets of argument values G_i , the action space \mathcal{A} is:

$$\mathcal{A} = L_{\text{func}} \times (L_{\text{arg1}} \times G_{\text{arg1}}) \times (L_{\text{arg2}} \times G_{\text{arg2}}) \times \cdots \times (L_{\text{argN}} \times G_{\text{argN}}) \quad (2)$$

where L_{func} is the set of function names on the page object in the Playwright API (Microsoft, 2024), and function arguments have a name and value $(L_{\text{arg1}} \times G_{\text{arg1}})$ corresponding to the Playwright API. We allow the agent access to call arbitrary functions in Playwright (Microsoft, 2024), a Microsoft-developed browser automation library that wraps a headless web browser. The agent’s goal is to complete a web navigation task specified via a natural language instruction $\mathbf{c} \in L$, starting from an initial URL, and operating the browser via function calls to the Playwright API until the desired task is complete, after which the agent calls the `stop` function and provides a final response:

$$\mathbf{a}_{\text{stop}} = (\text{“stop”}, (\text{“response”}, \text{“the task has been completed.”})) \quad (3)$$

We prompt the agent to produce a reasoning trace of a desired length (ablated in Figure 10) that contains function call actions as JSON in a fenced code block. To parse actions from the response, we employ a regex template that matches the first JSON code block, and a JSON decoder $f^{\text{text} \rightarrow \text{act}}(\cdot)$

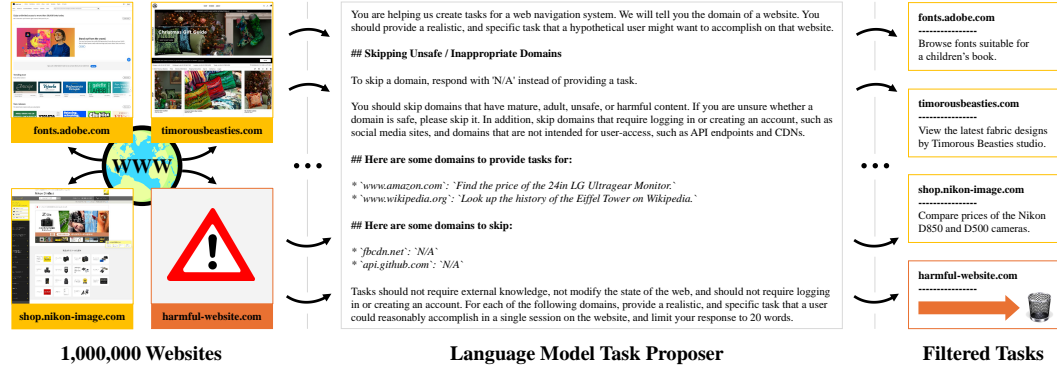


Figure 2: **Annotating 150k live sites with agentic tasks.** Starting from 1,000,000 websites, we employ a pretrained language model that marks sites as safe/unsafe for annotation, and assigns a realistic task that a hypothetical user might want to accomplish on each site. The task proposer aggressively filters out 85% of websites from the pipeline, resulting in 150k safe websites annotated with realistic tasks.

to parse the contents within the code block. When parsing fails due to invalid syntax, we allow the agent to generate a second response. Equipped with a language model agent that makes function calls with the Playwright API, we may consider the crucial task of obtaining large and diverse data.

4 INTERNET-SCALE TASK GENERATION

Training the next generation of LLM agents requires a large and diverse set of websites and tasks beyond what researchers have annotated so far (Deng et al., 2023; Zhou et al., 2024b; Koh et al., 2024a; Liu et al., 2024; Lù et al., 2024; Rawles et al., 2023; He et al., 2024). We develop an approach to efficiently annotate vast numbers of sites from diverse sections of the internet with agentic tasks. Our approach introduces two important desiderata: (1) it should not rely on human annotations, and (2) tasks should derive from a feedback process that deeply explores the environment.

4.1 LANGUAGE MODEL TASK PROPOSER

The key idea in stage one of the pipeline is a feedback loop, where a language model task proposer $\psi_{\text{LLM}}(\cdot)$ guides exploration on a website via an initial easy task. We then run a language model agent to explore the site, conditioned on the initial task, which produces an exploratory trajectory that is fed back to the task proposer. Conditioned on a trajectory that deeply explores the website, the task proposer then creates a harder, grounded task. This process is summarized as an equation.

$$\mathbf{c}_{n+1} \sim \psi_{\text{LLM}}(\cdot | \mathbf{w}, \mathbf{c}_n, \underbrace{\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T, \mathbf{r}_T}_{\text{last trajectory}}) \text{ st } \mathbf{a}_t \sim \pi_{\text{LLM}}(\cdot | \mathbf{s}_t, \mathbf{c}_n) \text{ s}_{t+1} \sim \mathbb{P}(\cdot | \mathbf{s}_t, \mathbf{a}_t) \quad (4)$$

The website url is \mathbf{w} , the initial task is \mathbf{c}_0 , the trajectory includes states \mathbf{s}_t , actions \mathbf{a}_t , judge score \mathbf{r}_T , and a harder, grounded task $\mathbf{c}_{n>0}$ for the next loop. Highlighted in Figure 2, we annotate 150k sites with tasks by scaling Equation 4, and release them on huggingface at: *data-for-agents/insta-150k-v2*.

Model Details. We utilize pretrained and frozen language models that conform to a chat interface and accept system, user, and assistant prompts. The task proposer system prompt is listed in Appendix 11, and details all cases for which sites are considered unsafe. We employ the Llama 3.1 family of LLMs from Meta (Grattafiori et al., 2024; Touvron et al., 2023b;a), the GPT family of LLMs from OpenAI, and the Gemini family of LLMs from Google. Inference is served using vLLM (Kwon et al., 2023) for the Llama series of models. We employ a sampling temperature of 0.5 and a maximum budget of 1024 new generated tokens, while all other parameters are kept as defaults in the OpenAI chat completions API, which is used to make inference calls to all LLMs.

Prompt Details. The task proposer operates in two phases. In an initial phase when just the url of a website is observed, the task proposer generates an initial task $\mathbf{c}_0 \sim \psi_{\text{LLM}}(\mathbf{c}_0 | \mathbf{w})$, and can mark a website as unsafe by setting $\mathbf{c}_0 = \text{N/A}$. The system prompt for this phase is listed in

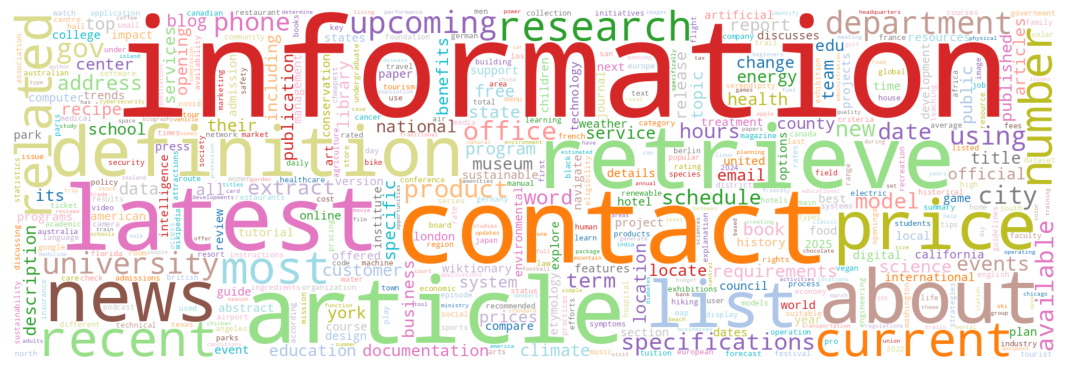


Figure 3: **Most frequent words in our tasks.** This wordcloud shows the top 500 most frequent words in tasks from the training set of our official huggingface dataset. The size of each word corresponds to its frequency in the dataset. Our tasks span diverse categories and lexicon.

Appendix 11. Agents discussed in Section 5 explore 150k sites annotated with initial tasks \mathbf{c}_0 and produce trajectories. In a second phase of task generation, we prompt the task proposer with the website url \mathbf{w} , the initial task \mathbf{c}_0 , the trajectory $\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T$, and a system prompt that instructs the LLM to produce a reasoning trace that contains a harder, grounded task $\mathbf{c}_{n>0}$. The system prompt for the second phase of task generation is listed in Appendix 11. The refined tasks produced by this iterative process lead to broadly capable agents, demonstrated in Section 6 by our ability to zero-shot transfer agents trained on our data to Web Voyager (He et al., 2024) and compete with frontier LLMs.

The design of the task proposer as a feedback process is important, and the full potential of this design will be realized in future work that trains agents with an on-policy reinforcement learning algorithm. In such future work, the task proposer can be used within the RL loop to generate incrementally harder tasks as agents learn. For this paper, we employ one loop of task generation.

4.2 SAFETY & RELIABILITY

Safety is critical when building autonomous agents. The internet contains significant amounts of content that should be removed from training data, in order to avoid agents learning harmful behaviors. To understand the robustness of our safety filter, we input 100 carefully selected websites to the task proposer, of which 50 contain harmful, or mature content. Table 1 reports the accuracy, precision, and recall of the safety filter on this data. For a variety of LLMs, the safety filter displays high accuracy—up to 97% of websites are correctly classified, and recall for detecting unsafe websites is as high as 1.0, suggesting that nearly all unsafe websites are detected by the safety filter.

Method	Acc.	Prec.	Recall
<i>Llama 3.1 70B</i>	85%	0.77	1.00
<i>GPT-4o</i>	95%	0.91	1.00
<i>Gemini 1.5 Pro</i>	97%	0.96	0.98

Table 1: The safety filter displays a high accuracy. We measure the accuracy, precision, and recall of the safety filter on a set of 100 websites, where 50 contain harmful, or mature content. Up to 97% of websites are correctly classified, and recall is as high as 1.0.

Method	Verifiable Rate
<i>Llama 3.1 70B</i>	75%
<i>GPT-4o</i>	85%
<i>Gemini 1.5 Pro</i>	89%

Table 2: Generated tasks are typically achievable. We measure the rate that human workers were able to complete and verify their completion of tasks produced by the task proposer for a set of 100 websites. Up to 89% of tasks are achievable, and verifiable.

Reliability is equally important for autonomous agents. Instructions should be followed faithfully by agents, and this requires training them with tasks that are achievable, and verifiable. Table 2 reports the rate that human workers were able to complete and verify their completion of tasks produced by the task proposer in its initial phase. Up to 89% of tasks are achievable, and verifiable according to the study, which suggests the pipeline is producing reliable tasks. Together with results in Section 6, it is likely that data from the InSTA pipeline leads to agents that follow instructions faithfully.

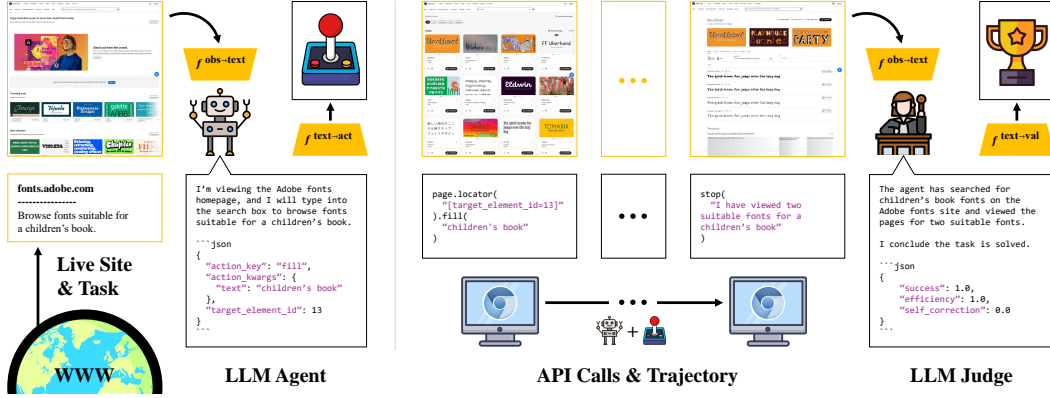


Figure 4: **Automatic evaluation for agents with language model judges.** Building on the large and diverse set of tasks generated by the pipeline, we employ pretrained language models to attempt and evaluate web navigation tasks. We dispatch language model agents to perform tasks by making calls to the Playwright API. We then employ language model judges to evaluate the trajectories.

4.3 SCALING TO 150,000 WEBSITES

We leverage Common Crawl for task generation. As of May 2025, the latest web graph released by The Common Crawl Foundation (2025) contains more than 300M unique hosts, which we adapt into a data source for agents. In particular, we select the top 1M sites based on their PageRank values. Common Crawl likely contain many unsafe websites, and these are filtered out by the task proposer. Each phase of task generation requires 14 hours of compute time using two 8-GPU v100 machines, and filters to 150k safe websites annotated with tasks. Statistics of tasks are shown in Figure 3.

5 INTERNET-SCALE ENVIRONMENT

By this point, we have reached our first goal—to improve coverage of real-world sites by annotating 150k diverse sites with challenging agentic tasks. To reach our second goal, and move beyond a static dataset, towards a *dynamic internet-scale environment*, we require a robust evaluator for these tasks. Evaluation presents a subtle challenge. The web evolves constantly, and daily changes in website content may invalidate a fixed ground truth reference solution. Driven by necessity, this environment must be evaluated by a model that judges whether an agent’s solution is correct, in the latest context.

5.1 EVALUATION WITH LANGUAGE MODELS

We model the process of evaluating trajectories from agents as a classification problem, where the goal is to estimate the probability r_T that a task c_n is solved, and generate r_T via next-token prediction, conditioned on a system prompt $\mathbf{x}_{\text{judge}}$, a task c_n , and a trajectory $s_1, \mathbf{a}_1, \dots, s_T, \mathbf{a}_T$. The LLM is instructed to produce a reasoning trace that scores the agent on a scale from 0 to 1, and embed the score as JSON in a fenced code block. We employ a regex template that matches to the first code block in the response, and a JSON decoder to parse r_T from the response, given by $f^{\text{text} \rightarrow \text{val}}(\cdot)$.

$$\mathbf{r}_T = f^{\text{text} \rightarrow \text{val}}(\text{LLM}(\mathbf{x}_{\text{judge}}, c_n, \text{Enc}(s_1), \mathbf{a}_1, \dots, \text{Enc}(s_T), \mathbf{a}_T)) \quad (5)$$

Verifying The Judge. To understand the robustness of the judge, we measure its accuracy detecting successful trajectories that were annotated by human workers. We annotate 100 trajectories with binary success labels, and apply a threshold $r_T > 0.5$ to obtain binary predictions from the judge. Figure 5 reports the accuracy of the judge as a function of the PageRank values of sites, and as a function of the confidence of the judge, given by $\text{conf} = 2 \cdot |r_T - 1/2|$. For all LLMs tested, the judge shows a high accuracy, ranging from 78.0% for *Gemini 1.5 Pro*, to 81.7% for *Llama 3.1 70B*, and 82.6% for *GPT-4o*. Accuracy is stable as PageRank falls, suggesting the judge is accurate for less popular sites that may not be well represented in the LLM’s training data. Shockingly, the confidence predicted by LLMs is highly interpretable, and correlates with accuracy to the point that trajectories where $\text{conf} = 1$ are classified with an accuracy up to 93.1%. The emergent robustness of the judge equips us to efficiently and accurately verify agent solutions on a dynamic internet.

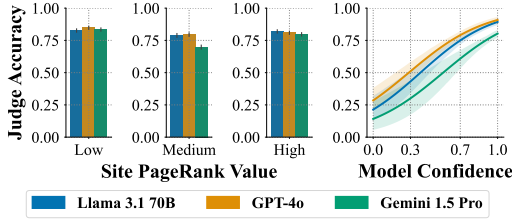


Figure 5: **Language models are robust evaluators.** We measure the accuracy of language models for detecting successful trajectories, and find that accuracy remains stable relative to PageRank values (*left plot*). As models become more confident, their accuracy improves (*right plot*), suggesting confidence is a useful proxy for the reliability of their predictions.

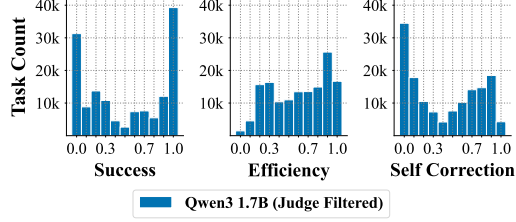


Figure 6: **Statistics for our agent reasoning dataset.** We conduct a large data collection experiment using our top checkpoint for *Qwen3 1.7B*. Our dataset has 2.2M screenshots, 2.2M reasoning traces for actions, and 150k traces for the judge. 50.0% of the trajectories are successful according to the judge, and have diverse ratings for efficiency and self-correction.

5.2 SCALING TO 150,000 AGENTS

With our task proposer, agent, and judge driven by pretrained language models, we have all components needed to harness internet-scale data. We conduct a large data collection experiment, running language model agents to complete tasks on 150k websites from our official dataset, producing 2.2M screenshots and 2.2M reasoning traces for actions within 150k trajectories. The judge annotates these trajectories, producing 150k reasoning traces for evaluations, and leading to the statistics in Figure 6. For this experiment, we employ a fine-tuned *Qwen3 1.7B* as the agent (refer to the next section), and *Qwen3 235B* zero-shot as the judge. Data collection requires 1,200 v100 GPU hours, and costs \$521.55 based on current AWS spot instance pricing, a fraction of the budget that industry labs are spending towards agents. If you have the right data, this small budget is sufficient and no human annotations are required to build models that compete with frontier LLMs as agents.

6 TRAINING AGENTS

We’ve built an internet-scale data flywheel that produces trajectories annotated with scores from a judge that can help us train LLM agents. To understand the quality of data produced by the flywheel, we conduct a series of experiments training models on the data, and testing on popular benchmarks. These experiments focus on three main questions: (1) *what is the impact of increasing data scale?* (2) *do agents transfer to new domains?* (3) *do agents scale with test-time compute?*

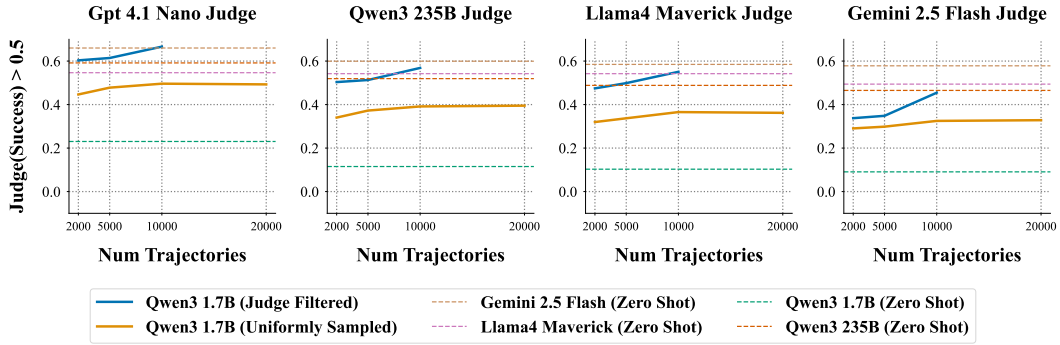


Figure 7: **InSTA unlocks great potential in small models.** We train agents based on *Qwen 3 1.7B* using trajectories produced by a *Qwen 3 235B* data collection policy, and optionally filtered by a *Qwen 3 235B* judge (see Judge Filtered vs. Uniformly Sampled). We report success rates on a test set of 3,000 held-out websites and tasks. Before training, *Qwen 3 1.7B* has a zero-shot success rate of 11.5% according to a *Qwen 3 235B* judge, and we improve this by +45.3% absolute percentage points. Our top checkpoint outperforms the *Qwen 3 235B* data collection policy, and *Llama 4 Maverick*, a frontier LLM with 400B parameters, for which our model is 235 times smaller. Notably, filtering with a *Qwen 3 235B* judge leads to agents that improve according to independent secondary judges, including *Gemini 2.5 Flash*, *Llama 4 Maverick*, and *Gpt 4.1 Nano*, suggesting it generalizes well.

6.1 PERFORMANCE IMPROVES WITH DATA SCALE

The recurring lesson in deep learning is that large-scale high-quality data wins, but researchers are struggling to materialize this promise for agents (Xue et al., 2025). Our paper aims to solve the data problem blocking researchers from materializing this promise, and this experiment provides a valuable signal that scaling high-quality data allows small models to compete with strong LLMs from top industry labs. To proceed, we collect 20k trajectories using a *Qwen 3 235B* data collection policy, annotated with scores from a *Qwen 3 235B* judge. We then train agents based on *Qwen 3 1.7B* with SFT on varying scales of the data. Results in Figure 7 show that performance improves with increasing data scale, and gains scale faster on data filtered by the judge. To filter the data, we select trajectories where $\text{Judge}(\text{Success}) = 1$. Our top checkpoint outperforms the *Qwen 3 235B* data collection policy, and beats *Llama 4 Maverick*, a frontier LLM with 400B parameters, for which our model is 235 times smaller. The trend in the figure suggests there is room to scale further.

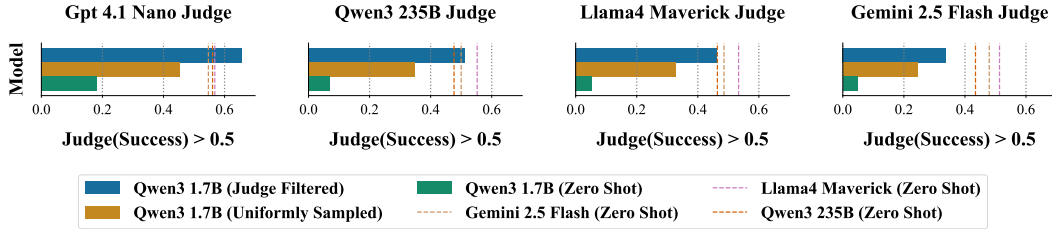


Figure 8: **Our agents zero-shot transfer to WebVoyager.** With no additional training or specialized data, our checkpoints for *Qwen 3 1.7B* in Section 6.1 zero-shot transfer to the WebVoyager benchmark. Trends found on our test set appear to hold for WebVoyager as well, and our top checkpoints for *Qwen 3 1.7B* continue to match frontier LLMs in performance for three of four judges.

6.2 AGENTS TRANSFER TO NEW DOMAINS

Statistically correct evaluation for deep learning models requires a test dataset that does not overlap with the training dataset, but researchers do not agree on how to implement this guideline for agents. Recent works train agents on the same websites they test on (Murty et al., 2025; Zhou et al., 2024c; Putta et al., 2024), which may obfuscate progress (Xue et al., 2025). Our next experiment shows that we can implement a stronger train-test split. Agents trained with our data can zero-shot transfer to WebVoyager (He et al., 2024) without any data from the benchmark. Results in Figure 8 show our *Qwen 3 1.7B* checkpoint matching strong LLMs on WebVoyager (He et al., 2024) for three of four judges, confirming trends in Section 6.1. Our ability to zero-shot transfer relatively small agents to WebVoyager (He et al., 2024) makes it likely that our pipeline leads to capable agents.

Static Benchmarks. To complement the online evaluation in Figure 8, we also explore how our data impacts agents trained on static benchmarks. We first train baseline agents on human demonstrations from the Mind2Web (Deng et al., 2023) and WebLINX (Lù et al., 2024) datasets. We write a preprocessor to convert our data into the expected action-observation formats these benchmarks use (which involves discarding our reasoning trace). With data converted, we compare the baseline to agents trained on a mix of 80% human data, and 20% our data, and test on (1) their official test set, and (2) 500 diverse sites from our official test set. Results in Figure 9 show that agents trained with our data perform equally well on the original test sets for these benchmark, but generalize better to our harder test set. Overall, we see +149.0% for WebLINX agents, +156.3% for Mind2Web agents, and gains in *Step Accuracy* on our test set are larger for the harder tasks. Additional experimental details are listed in Appendix H. On three popular benchmarks (WebVoyager, WebLINX, Mind2Web), our pipeline trains capable agents, and does not overly rely on human annotations.

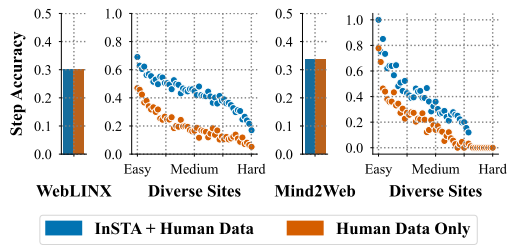


Figure 9: **Our data transfers to static benchmarks.** We train agents with all human data from the WebLINX and Mind2Web training sets, and resulting agents struggle to generalize to more diverse test data. Adding our data improves generalization by +149.0% for WebLINX, and +156.3% for Mind2Web.

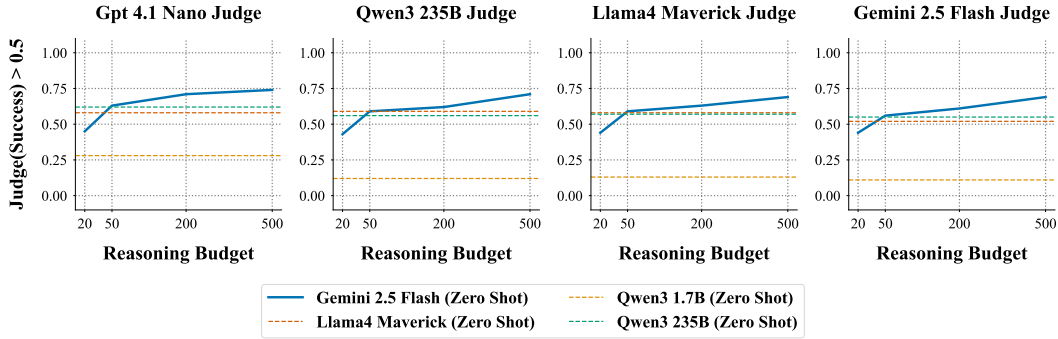


Figure 10: **Agents improve with a larger reasoning budget.** We ablate the number of tokens in the reasoning budget for the top-performing agent, and see a monotonic improvement in the success rate as the reasoning budget increases. *Gemini 2.5 Flash* has a 70% success rate with a budget of 500 reasoning tokens, up from 60% for a budget of 50 tokens. The scaling of performance with the reasoning budget highlights a promising behavior in successful web agents that can be studied.

6.3 PERFORMANCE SCALES WITH TEST-TIME COMPUTE

To understand the ability for web agents to scale with additional test-time compute, we ablate the number of tokens in the reasoning budget. Figure 10 shows the success rate as a function of the reasoning budget for *Gemini 2.5 Flash*, the top-performing agent we tested. There is a monotonic improvement in the success rate as the reasoning budget increases, and the trend suggests that performance may not be saturated with 500 reasoning tokens. Training language model agents to reason before taking actions is a promising path to better agents, and we are releasing a large reasoning dataset for multimodal agents to study this. Our dataset contains 2.2M screenshots, 2.2M reasoning traces for actions, 150k traces for judge evaluations, and led to the results in Section 6.1. The data will be linked on our website, alongside an official huggingface dataset for tasks.

7 CONCLUSION

In the spirit of deep learning, we have developed an approach to efficiently harness internet data for LLM agents, and have unlocked a *dynamic internet-scale environment*. In building this environment, we presented a method to annotate 150k diverse sites with challenging agentic tasks, and showed how training on data from our pipeline allows small models to compete with frontier LLMs as agents, on a fraction of the budget. Our pipeline consists of a task proposer, agent, and judge driven by pretrained language models that together curate high-quality data for agents, without human intervention.

Our top checkpoint for *Qwen 3 1.7B* has a success rate of 56.9% on our test environment, outperforming the data collection policy *Qwen 3 235B*, beating the 235 times larger *Llama 4 Maverick*, and reaching 94.7% of the performance of *Gemini 2.5 Flash*, while being smaller and faster than these. Our models zero-shot transfer to WebVoyager, and scale with test-time compute. We are releasing the entire pipeline, including code, models and data, so that it may serve as a foundation for researchers to build the next generation of language model agents with internet data.

7.1 FUTURE WORK

Our work reveals several exciting directions in future work. First, our work can be scaled further: the latest Common Crawl release contains data for more than 300M sites, suggesting another 1,000 times more data could be available for agents by scaling the pipeline. In addition, we trained agents to optimize the judge scores indirectly via filtered SFT, and the judge’s high accuracy suggests that it could be optimized via reinforcement learning instead. RL is especially promising for how it can improve reasoning capabilities in agents. Finally, while the data we collect is multimodal, we focus on textual tasks in this paper, and our pipeline could be extended to produce multimodal tasks.

REFERENCES

- Jacob Andreas. Language models as agent models. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 5769–5779, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-emnlp.423. URL <https://aclanthology.org/2022.findings-emnlp.423>.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefer. Graph of thoughts: Solving elaborate problems with large language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17682–17690, Mar. 2024. doi: 10.1609/aaai.v38i16.29720. URL <https://ojs.aaai.org/index.php/AAAI/article/view/29720>.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrmke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. Sparks of artificial general intelligence: Early experiments with gpt-4, 2023. URL <https://arxiv.org/abs/2303.12712>.
- Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier, Karthik Narasimhan, and Shunyu Yao. Fireact: Toward language agent fine-tuning, 2023. URL <https://arxiv.org/abs/2310.05915>.
- Thibault Le Sellier De Chezelles, Maxime Gasse, Alexandre Drouin, Massimo Caccia, Léo Boisvert, Megh Thakkar, Tom Marty, Rim Assouel, Sahar Omid Shayegan, Lawrence Keunho Jang, Xing Han Lù, Ori Yoran, Dehan Kong, Frank F. Xu, Siva Reddy, Quentin Cappart, Graham Neubig, Ruslan Salakhutdinov, Nicolas Chapados, and Alexandre Lacoste. The browsergym ecosystem for web agent research, 2024. URL <https://arxiv.org/abs/2412.05467>.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web, 2023. URL <https://arxiv.org/abs/2306.06070>.
- Saumya Gandhi, Ritu Gala, Vijay Viswanathan, Tongshuang Wu, and Graham Neubig. Better synthetic data by retrieving and transforming existing datasets, 2024. URL <https://arxiv.org/abs/2404.14361>.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, and et al. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- 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. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6864–6890, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.371. URL <https://aclanthology.org/2024.acl-long.371>.
- Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent: A visual language model for gui agents, 2023. URL <https://arxiv.org/abs/2312.08914>.
- Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashmi Rungta, Krithika Iyer, Yuning Mao, Michael Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, and Madian Khabza. Llama guard: Llm-based input-output safeguard for human-ai conversations, 2023. URL <https://arxiv.org/abs/2312.06674>.

- 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, 2024a. URL <https://arxiv.org/abs/2401.13649>.
- Jing Yu Koh, Stephen McAleer, Daniel Fried, and Ruslan Salakhutdinov. Tree search for language model agents, 2024b. URL <https://arxiv.org/abs/2407.01476>.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention, 2023. URL <https://arxiv.org/abs/2309.06180>.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, and Sushant Prakash. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback, 2024. URL <https://arxiv.org/abs/2309.00267>.
- Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, Kai Shu, Lu Cheng, and Huan Liu. From generation to judgment: Opportunities and challenges of llm-as-a-judge. *arXiv preprint arXiv:2411.16594*, 2024.
- Junpeng Liu, Yifan Song, Bill Yuchen Lin, Wai Lam, Graham Neubig, Yuanzhi Li, and Xiang Yue. Visualwebbench: How far have multimodal llms evolved in web page understanding and grounding?, 2024. URL <https://arxiv.org/abs/2404.05955>.
- Xing Han Lù, Zdeněk Kasner, and Siva Reddy. Weblinx: Real-world website navigation with multi-turn dialogue, 2024. URL <https://arxiv.org/abs/2402.05930>.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=S37hOerQLB>.
- Microsoft. Playwright. <https://github.com/microsoft/playwright>, 2024.
- Arindam Mitra, Luciano Del Corro, Guoqing Zheng, Shweti Mahajan, Dany Rouhana, Andres Cotas, Yadong Lu, Wei ge Chen, Olga Vrousos, Corby Rosset, Fillipe Silva, Hamed Khanpour, Yash Lara, and Ahmed Awadallah. Agentinstruct: Toward generative teaching with agentic flows, 2024. URL <https://arxiv.org/abs/2407.03502>.
- Shikhar Murty, Hao Zhu, Dzmitry Bahdanau, and Christopher D. Manning. Nnetnav: Unsupervised learning of browser agents through environment interaction in the wild, 2025. URL <https://arxiv.org/abs/2410.02907>.
- Tianyue Ou, Frank F. Xu, Aman Madaan, Jiarui Liu, Robert Lo, Abishek Sridhar, Sudipta Sengupta, Dan Roth, Graham Neubig, and Shuyan Zhou. Synatra: Turning indirect knowledge into direct demonstrations for digital agents at scale, 2024. URL <https://arxiv.org/abs/2409.15637>.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS ’22*, Red Hook, NY, USA, 2024. Curran Associates Inc. ISBN 9781713871088.
- Ajay Patel, Markus Hofmarcher, Claudiu Leoveanu-Condrei, Marius-Constantin Dinu, Chris Callison-Burch, and Sepp Hochreiter. Large language models can self-improve at web agent tasks, 2024. URL <https://arxiv.org/abs/2405.20309>.

- Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, and Boi Faltings. REFINER: Reasoning feedback on intermediate representations. In Yvette Graham and Matthew Purver (eds.), *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1100–1126, St. Julian’s, Malta, March 2024. Association for Computational Linguistics. URL <https://aclanthology.org/2024.eacl-long.67>.
- Pranav Putta, Edmund Mills, Naman Garg, Sumeet Motwani, Chelsea Finn, Divyansh Garg, and Rafael Rafailov. Agent q: Advanced reasoning and learning for autonomous ai agents, 2024. URL <https://arxiv.org/abs/2408.07199>.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners, 2019. URL <https://api.semanticscholar.org/CorpusID:160025533>.
- Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. Android in the wild: A large-scale dataset for android device control, 2023. URL <https://arxiv.org/abs/2307.10088>.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=Yacmpz84TH>.
- Amrith Setlur, Saurabh Garg, Xinyang Geng, Naman Garg, Virginia Smith, and Aviral Kumar. RL on incorrect synthetic data scales the efficiency of llm math reasoning by eight-fold, 2024. URL <https://arxiv.org/abs/2406.14532>.
- Junhong Shen, Atishay Jain, Zedian Xiao, Ishan Amlekar, Mouad Hadji, Aaron Podolny, and Ameet Talwalkar. Scribeagent: Towards specialized web agents using production-scale workflow data, 2024. URL <https://arxiv.org/abs/2411.15004>.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters, 2024. URL <https://arxiv.org/abs/2408.03314>.
- Hanshi Sun, Momin Haider, Ruiqi Zhang, Huitao Yang, Jiahao Qiu, Ming Yin, Mengdi Wang, Peter Bartlett, and Andrea Zanette. Fast best-of-n decoding via speculative rejection. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=348hfcprUs>.
- Fahim Tajwar, Anikait Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie, Stefano Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of llms should leverage suboptimal, on-policy data, 2024. URL <https://arxiv.org/abs/2404.14367>.
- The Common Crawl Foundation. Common crawl, 2025. URL <https://commoncrawl.org/>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023a. URL <https://arxiv.org/abs/2302.13971>.
- Hugo Touvron, Louis Martin, Kevin Stone, and et al. Llama 2: Open foundation and fine-tuned chat models, 2023b. URL <https://arxiv.org/abs/2307.09288>.
- Brandon Trabucco, Kyle Doherty, Max A Gurinas, and Ruslan Salakhutdinov. Effective data augmentation with diffusion models. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=ZWzUA9zeAg>.
- Karthik Valmeekam, Kaya Stechly, and Subbarao Kambhampati. Llms still can’t plan; can lrms? a preliminary evaluation of openai’s o1 on planbench, 2024. URL <https://arxiv.org/abs/2409.13373>.

- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Jirong Wen. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6), March 2024. ISSN 2095-2236. doi: 10.1007/s11704-024-40231-1. URL <http://dx.doi.org/10.1007/s11704-024-40231-1>.
- Junlin Xie, Zhihong Chen, Ruifei Zhang, Xiang Wan, and Guanbin Li. Large multimodal agents: A survey, 2024. URL <https://arxiv.org/abs/2402.15116>.
- 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. 2025. URL <https://arxiv.org/abs/2504.01382>.
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents, 2023a. URL <https://arxiv.org/abs/2207.01206>.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models, 2023b. URL <https://arxiv.org/abs/2305.10601>.
- Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Zhi Huang, Carlos Guestrin, and James Zou. Textgrad: Automatic "differentiation" via text, 2024. URL <https://arxiv.org/abs/2406.07496>.
- Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. Agenttuning: Enabling generalized agent abilities for llms, 2023. URL <https://arxiv.org/abs/2310.12823>.
- Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent: Multimodal agents as smartphone users, 2023. URL <https://arxiv.org/abs/2312.13771>.
- Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. Generative verifiers: Reward modeling as next-token prediction, 2024. URL <https://arxiv.org/abs/2408.15240>.
- Tianyang Zhong, Zhengliang Liu, Yi Pan, and et al. Evaluation of openai o1: Opportunities and challenges of agi, 2024. URL <https://arxiv.org/abs/2409.18486>.
- Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Language agent tree search unifies reasoning acting and planning in language models, 2024a. URL <https://arxiv.org/abs/2310.04406>.
- Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web environment for building autonomous agents, 2024b. URL <https://arxiv.org/abs/2307.13854>.
- Yifei Zhou, Qianlan Yang, Kaixiang Lin, Min Bai, Xiong Zhou, Yu-Xiong Wang, Sergey Levine, and Erran Li. Proposer-agent-evaluator(pae): Autonomous skill discovery for foundation model internet agents, 2024c. URL <https://arxiv.org/abs/2412.13194>.

A LIMITATIONS & SAFEGUARDS

Language model agents present unique challenges and risks when applied to live tasks on the internet. For instance, agents visiting shopping sites can influence the statistics produced by analytics tools, which can impact prices on products, and product decisions from companies. Furthermore, agents seeing harmful content on the web can add that content to datasets inadvertently, and propagate harmful behaviors to future agents. We mitigate these risks in the design of the task proposal stage. We consider the risks posed to analytics tools by limiting the engagement between agents and sites. We generate only one task per website, and we limit agents to just 30 actions per site, which includes clicks, typing, dropdown selection, and more. By limiting the interaction between agents and sites, the change in website traffic generated by the InSTA pipeline is minimal (just 90 seconds of interaction per site on average). By utilizing data from the InSTA pipeline in an offline fashion, as in Section 6 of the main paper, no additional web traffic is generated when training agents. To ensure that agents do not modify the state of the web (i.e. avoid attempting to make purchases, avoid leaving comments on posts, avoid making accounts, etc), we provide instruct the task proposer (see Figure 2) to avoid writing tasks that require the agent to interact with personal data, or user accounts.

The task proposer is instructed via the system prompt to filter out sites with harmful content, sites not intended for user access, and sites that require making an account to operate, including social media, and forum sites. There is likely a manner to safely train agents to operate user accounts, but we leave this task to future researchers. We explore the performance of the task proposer at filtering out unsuitable sites in Section 4.2, and find that all models detect unsuitable sites with a recall from 0.98 to 1.0, and accuracy up to 97%, suggesting our filter is reliable. Sites used to benchmark the performance of the safety filter are discussed in Appendix E, and thoroughly test the safety filter.

To remove Personally Identifiable Information (PII) from the data used for training agents, we include `scrubadub`, an industry standard PII removal tool for python developed by Leap Beyond, a data consultancy based in the European Union. Our pipeline has an argument that toggles the usage of `scrubadub` to remove PII from all website data, and we recommend this option be set.

B ETHICAL CONSIDERATIONS

One important ethical consideration when harnessing internet data is to carefully handle copyrighted, private, and sensitive materials. The internet contains vast amounts of PII, which should be avoided when training models. We address this in two ways. First, we instruct the task proposer to filter out sites that may contain PII, including social media websites, and forums. Second, our pipeline has an argument that toggles the usage of `scrubadub` to remove PII, and we recommend this option be set. To prevent copyrighted materials from being used for training, we do not allow the agent to copy any books, documents, videos, audio, or files. The information present in the HTML DOM is publically available, and is converted into a sparse Markdown summary that focuses on interactive elements, and significantly transforms the underlying webpage in a manner that is not reversible.

These steps significantly reduce, but do not completely eliminate the risk that private, and sensitive materials are shown to the agent, and methods for detecting, replacing, and removing such materials from text and images remains an important task for researchers working on safety.

C BROADER IMPACTS

As their capabilities broaden, language models are being used to interface with real-world systems. This shift comes with several benefits and risks. Agents that operate your computer to aid in work tasks can significantly boost productivity for certain workers, but can displace others whose jobs have been fully automated. Agents that operate browsers to complete personal tasks provide convenience, but expose a new attack vector where rogue agents perform unintended actions. Certain risks can be mitigated with proper safeguards, such as post-processing data to prevent jail-breaking, but other risks are existential, harder to address purely with new agent research, and may require policy changes.

```

You are helping us create tasks for a web navigation system. We
  ↳ will tell you the domain of a website. You should provide a
  ↳ realistic, and specific task that a hypothetical user might
  ↳ want to accomplish on that website.

## Skipping Unsafe / Inappropriate Domains

To skip a domain, respond with 'N/A' instead of providing a task.

You should skip domains that have mature, adult, unsafe, or
  ↳ harmful content. If you are unsure whether a domain is safe,
  ↳ please skip it. In addition, skip domains that require
  ↳ logging in or creating an account, such as social media
  ↳ sites, and domains that are not intended for user-access,
  ↳ such as API endpoints and CDNs.

## Here are some domains to provide tasks for:

* 'www.amazon.com': 'Find the price of the 24in LG Ultragear
  ↳ Monitor.'
* 'www.wikipedia.org': 'Look up the history of the Eiffel Tower on
  ↳ Wikipedia.'

## Here are some domains to skip:

* 'fbcdn.net': 'N/A'
* 'api.github.com': 'N/A'

Tasks should not require external knowledge, not modify the state
  ↳ of the web, and should not require logging in or creating an
  ↳ account. For each of the following domains, provide a
  ↳ realistic, and specific task that a user could reasonably
  ↳ accomplish in a single session on the website, and limit
  ↳ your response to 20 words.

```

Figure 11: **System prompt for the exploration phase of task generation.** We design the system prompt for task generation to detect and remove unsafe websites. This prompt ensures that tasks are passive, and do not modify content on a website. Refer to the next figures for the in-context examples used for the task proposer, and the system prompt used in the feedback step.

D AGENTS.TXT & STANDARDS FOR INTERNET AGENTS

Akin to `robots.txt` directives, website creators should have a standard format to specify how internet agents are allowed to interact with their websites, and what information on webpages agents are allowed to see. Desirable controls include rate limits for interactions, limits for maximum numbers of interactions, restrictions to allow agents to interact with certain pages and not others, and restrictions on the kind of data on webpages that agents are allowed to observe (achieved via tagging elements to hide their content from agents). In addition to restricting the data available to agents, website creators should have the ability to specify locations for “playgrounds” that replicate certain key functions of their site with virtual tasks and simulated data that are intended to teach agents how to operate their site while directing traffic from agents away from user-facing pages.

E MORE DETAILS ON TASK GENERATION

We provide the system prompt used in the first phase of the task generation loop in Figure 11. This prompt was provided to Llama 3.1 70B, GPT-4o, and Gemini 1.5 Pro to generate tasks and filter

sites unsuitable for annotation in Section 4. We carefully designed this system prompt to enforce that generated tasks are passive, and do not modify content on a website. In addition to this system prompt, we employed a list of 100 hand-picked in-context examples of website URLs and appropriate tasks, which are provided in the following JSON list. When querying an LLM, we randomly sample 16 in-context examples from the list, and provide only these examples to the LLM to generate a task to guide exploration of the site. This improves diversity in the exploration phase.

```
[
  {
    "domain": "archive.org",
    "task": "Identify the oldest book available in the public
    ↪ domain on this site."
  },
  {
    "domain": "arxiv.org",
    "task": "Retrieve the latest preprint paper on machine
    ↪ learning."
  },
  {
    "domain": "wikibooks.org",
    "task": "Find a freely available textbook on linear algebra
    ↪ ."
  },
  {
    "domain": "wiktionary.org",
    "task": "Get the definition and etymology of the word '
    ↪ serendipity'."
  },
  {
    "domain": "openlibrary.org",
    "task": "Locate an ebook about classic literature that is
    ↪ available for borrowing."
  },
  {
    "domain": "openculture.com",
    "task": "Find a free online course on ancient history."
  },
  {
    "domain": "theguardian.com",
    "task": "Retrieve an article discussing recent trends in
    ↪ renewable energy."
  },
  {
    "domain": "medium.com",
    "task": "Identify a highly rated blog post on productivity
    ↪ hacks."
  },
  {
    "domain": "goodreads.com",
    "task": "Find the most popular book related to neuroscience
    ↪ ."
  },
  {
    "domain": "wired.com",
    "task": "Retrieve an article about the latest advancements
    ↪ in wearable technology."
  },
]
```

```

864     "domain": "data.gov",
865     "task": "Identify the latest government dataset on climate
866           ↪ change."
867 },
868 {
869     "domain": "kaggle.com",
870     "task": "Find a well-documented data science competition on
871           ↪ image recognition."
872 },
873 {
874     "domain": "gov.uk",
875     "task": "Locate the latest UK government report on
876           ↪ healthcare."
877 },
878 {
879     "domain": "unsplash.com",
880     "task": "Find a high-resolution image of the Milky Way
881           ↪ Galaxy."
882 },
883 {
884     "domain": "pexels.com",
885     "task": "Retrieve a popular photo tagged with 'nature'."
886 },
887 {
888     "domain": "creativecommons.org",
889     "task": "Find an article explaining Creative Commons
890           ↪ licensing types."
891 },
892 {
893     "domain": "pypi.org",
894     "task": "Retrieve the most downloaded Python package for
895           ↪ data analysis."
896 },
897 {
898     "domain": "huggingface.co",
899     "task": "Identify a popular machine learning model on this
900           ↪ platform."
901 },
902 {
903     "domain": "sciencenews.org",
904     "task": "Find the most recent article on the health impacts
905           ↪ of air pollution."
906 },
907 {
908     "domain": "mit.edu",
909     "task": "Retrieve a publicly available research paper on
910           ↪ quantum computing."
911 },
912 {
913     "domain": "springer.com",
914     "task": "Identify the latest edition of a Springer book on
915           ↪ robotics."
916 },
917 {
918     "domain": "jstor.org",
919     "task": "Find a research paper discussing the history of the
920           ↪ Internet."
921 },
922 {

```

```

918     "domain": "biorxiv.org",
919     "task": "Retrieve the most recent bioRxiv preprint on CRISPR
920           ↪ technology."
921   },
922   {
923     "domain": "medrxiv.org",
924     "task": "Find a public health preprint related to COVID-19."
925   },
926   {
927     "domain": "commons.wikimedia.org",
928     "task": "Retrieve a high-resolution image of the Eiffel
929           ↪ Tower."
930   },
931   {
932     "domain": "scholar.google.com",
933     "task": "Find the most cited article by a specific
934           ↪ researcher."
935   },
936   {
937     "domain": "plos.org",
938     "task": "Locate the latest research paper on gene editing
939           ↪ published here."
940   },
941   {
942     "domain": "flickr.com",
943     "task": "Find a photo that has been released under a
944           ↪ Creative Commons license."
945   },
946   {
947     "domain": "datacite.org",
948     "task": "Retrieve metadata for a dataset related to
949           ↪ environmental studies."
950   },
951   {
952     "domain": "orcid.org",
953     "task": "Find the ORCID ID of a well-known researcher in AI
954           ↪ ."
955   },
956   {
957     "domain": "zotero.org",
958     "task": "Retrieve an article discussing citation management
959           ↪ tools."
960   },
961   {
962     "domain": "github.com",
963     "task": "Find the most starred repository on deep learning."
964   },
965   {
966     "domain": "figshare.com",
967     "task": "Retrieve an open dataset on climate patterns."
968   },
969   {
970     "domain": "zenodo.org",
971     "task": "Find the latest publication on open science
           ↪ practices."

```



```

972 },
973 {
974   "domain": "biodiversitylibrary.org",
975   "task": "Retrieve a scanned copy of an 18th-century
976           ↪ botanical illustration."
977 },
978 {
979   "domain": "genome.gov",
980   "task": "Find the latest update on the Human Genome Project
981           ↪ ."
982 },
983 {
984   "domain": "merriam-webster.com",
985   "task": "Retrieve the definition and usage of the word '
986           ↪ quantum'."
987 },
988 {
989   "domain": "stanford.edu",
990   "task": "Find the most recent online lecture on artificial
991           ↪ intelligence."
992 },
993 {
994   "domain": "edx.org",
995   "task": "Retrieve a TED Talk on leadership in technology."
996 },
997 {
998   "domain": "ted.com",
999   "task": "Find the latest ocean temperature data available."
1000 },
1001 {
1002   "domain": "noaa.gov",
1003   "task": "Retrieve a dataset related to consumer behavior."
1004 },
1005 {
1006   "domain": "data.world",
1007   "task": "Find a course on data visualization."
1008 },
1009 {
1010   "domain": "curious.com",
1011   "task": "Retrieve a well-cited article on the psychological
1012           ↪ impact of social media."
1013 },
1014 {
1015   "domain": "theconversation.com",
1016   "task": "Identify a recent research paper on biodiversity
1017           ↪ conservation."
1018 },
1019 {
1020   "domain": "nature.com",
1021   "task": "Retrieve the latest article on genomics research."
1022 },
1023 {
1024   "domain": "pnas.org",
1025   "task": "Find a science news article on robotics
           ↪ advancements."
1026 },
1027 {
1028   "domain": "sciencedaily.com",
1029   "task": "Identify the top story on global health issues."

```

```

1026     },
1027     {
1028         "domain": "bbc.com",
1029         "task": "Retrieve a recent podcast episode about space
1030                 ↪ exploration."
1031     },
1032     {
1033         "domain": "npr.org",
1034         "task": "Locate the most recent update on the global
1035                 ↪ biodiversity status."
1036     }
1037 ]

```

We also provide the system prompt used in the second phase of the task generation loop, where trajectories from agents are fed back to the task proposer, which generates a harder, grounded task. This prompt instructs the task proposer to create a challenging task based on how an expert user could be expected to use the shown website. The task proposer also predicts a list on intermediate steps that can be used as a hint for agents, and a success criteria that can be used to improve the verifier.

```

1046 You are a helpful assistant designing tasks for a web automation
1047     ↪ script. I will show you previous runs of the script,
1048     ↪ including previous tasks, webpages, actions, and performance
1049     ↪ reviews, formatted in markdown. Help me design *challenging
1050     ↪ * new tasks.
1051
1052 ## Formatting The Proposed Task
1053
1054 Format your task in the following JSON schema:
1055
1056 ```json
1057 {
1058     "proposed_task": str,
1059     "steps": List[str],
1060     "criteria": str
1061 }
1062 ```
1063
1064 Here is what each key means:
1065
1066 - 'proposed_task': A specific, challenging task that an expert
1067     ↪ user might leverage this website to complete.
1068     - Must not require making an account, logging in, submitting
1069     ↪ personal information, making a purchase, or placing an
1070     ↪ order.
1071
1072 - 'steps': Steps an expert user would follow to complete the
1073     ↪ proposed task.
1074
1075 - 'criteria': The required answer, and criteria to determine if
1076     ↪ the task was completed.
1077
1078 ## Example Tasks For Inspiration
1079
1080 Suppose you want to design a task around the 'C-to-C Hose-Shut-Off
1081     ↪ Valve' on 'awg-fittings.com':
1082
1083 ```json
1084 {

```

```

1080     "proposed_task": "What is the C-to-C Hose-Shut-Off Valve length
1081         ↪ in mm?",
1082     "steps": [
1083         "Navigate to 'awg-fittings.com'",
1084         "Open the product catalog for fittings",
1085         "Locate the product listing for the C-to-C Hose-Shut-Off
1086         ↪ Valve",
1087         "Find the product length in mm, and respond with that length
1088         ↪ in the answer"
1089     ],
1090     "criteria": "The answer should include the specific length of
1091         ↪ '237 mm' for this product"
1092 }
1093 ...
1094 Suppose you want to design a task around the document 'The Angora
1095     ↪ cat; how to breed train and keep it' on 'biodiversitylibrary.
1096     ↪ org':
1097
1098 ```json
1099 {
1100     "proposed_task": "Open a scanned copy of 'The Angora cat; how
1101         ↪ to breed train and keep it'.",
1102     "steps": [
1103         "Navigate to 'biodiversitylibrary.org'",
1104         "Search for 'The Angora cat; how to breed train and keep it'
1105         ↪ in the search bar",
1106         "Click on the title of the document in the search results",
1107         "Confirm the correct document is displayed in an embedded
1108         ↪ PDF reader"
1109     ],
1110     "criteria": "The final webpage should display the correct
1111         ↪ document in an embedded PDF reader"
1112 }
1113 ...
1114 Suppose you want to design a task around the 'Generative
1115     ↪ Adversarial Networks' paper on 'scholar.google.com':
1116
1117 ```json
1118 {
1119     "proposed_task": "How many citations does the paper 'Generative
1120         ↪ Adversarial Networks' have?",
1121     "steps": [
1122         "Navigate to 'scholar.google.com'",
1123         "Search for 'Generative Adversarial Networks' in the search
1124         ↪ bar",
1125         "Locate the correct paper in the search results",
1126         "Find an up-to-date citation count, and respond with that
1127         ↪ count in the answer"
1128     ],
1129     "criteria": "The answer should include an up-to-date citation
1130         ↪ count, which is '80613' as of April 2025"
1131 }
1132 ...
1133 Suppose you want to design a task around the word 'serendipity' on
1134     ↪ 'wiktionary.org':

```

```

1134 ```json
1135 {
1136     "proposed_task": "What is the definition and etymology of the
1137     ↪ word 'serendipity'?",
1138     "steps": [
1139         "Navigate to 'wiktionary.org'",
1140         "Search for 'serendipity' in the search bar",
1141         "Find the definition and etymology sections of the '
1142         ↪ serendipity' page",
1143         "Summarize the contents of these sections in the answer"
1144     ],
1145     "criteria": "The answer should mention Serendip (or Serendib),
1146     ↪ coined by English writer and politician Horace Walpole in
1147     ↪ 1754"
1148 }
1149 ```
1150 Thanks for helping me design challenging new tasks, please follow
1151 ↪ the instructions carefully. Start your response with an
1152 ↪ analysis for how an expert user would leverage this website,
1153 ↪ followed by a step-by-step breakdown of your proposed task,
1154 ↪ and finally, enter your task in the JSON format. Respond in
1155 ↪ 500 words.

```

1157 E.1 DETAILS FOR SAFETY EXPERIMENT

1159 Using these prompts for task generation, we remove unsafe websites. To evaluate the performance
 1160 of our filter, we employed a set of 100 curated websites, where 50 are manually verified as safe,
 1161 and 50 are manually verified as unsafe based on the filtering conditions. These sites were chosen to
 1162 span popular sites that typical users are likely familiar with, and less popular websites that may be
 1163 underrepresented in LLM training data.

```

1165 safe_sites_list = ['dhss.mo.gov', 'dizionari.corriere.it', '
1166 ↪ southgippsland.vic.gov.au', 'ds.iris.edu', 'lobbycontrol.de
1167 ↪ ', '4rsmokehouse.com', 'barnsleyfc.co.uk', 'wiwi.uni-
1168 ↪ wuerzburg.de', 'uplandca.gov', 'lsus.edu', 'wpcode.com', '
1169 ↪ webopedia.internet.com', 'tamko.com', 'premierchristian.news
1170 ↪ ', 'genome.jgi.doe.gov', 'burgerking.ca', 'thehugoawards.org
1171 ↪ ', 'radio.fm', 'thevinyldistrict.com', 'unilang.org', '
1172 ↪ raywhitegroup.com', 'grapevinetexas.gov', 'sanfrancisco.
1173 ↪ cbslocal.com', 'hyde-design.co.uk', 'breastcancerfoundation.
1174 ↪ org.nz', 'ludwigsburg.de', 'ignitionrobotics.org', '
1175 ↪ deliverit.com.au', 'kodokan.org', 'clickstay.com', '
1176 ↪ searchdatamanagement.techtarget.com', 'oceanario.pt', '
1177 ↪ wentworthpuzzles.com', 'catholicworldreport.com', 'quizlet.
1178 ↪ com', 'innovation.nhs.uk', 'synonyms.reverso.net', 'news.
1179 ↪ siemens.co.uk', 'readability-score.com', 'co.modoc.ca.us', '
1180 ↪ cityofmyrtlebeach.com', 'loire.gouv.fr', 'lawphil.net', '
1181 ↪ saem.org', 'parmigianoreggiano.it', 'engaging-data.com', '
1182 ↪ itf-tkd.org', 'aka.education.gov.uk', 'ub.uni-kl.de', '
1183 ↪ mottchildren.org']
1184 unsafe_sites_list = ['refer to the code release']

```

1186 We provide the list of unsafe sites in the code to keep this manuscript clean. Using these lists, we
 1187 conduct an experiment measuring the accuracy, precision, and recall of our safety filter for detecting
 sites that are not suitable for training agents. In this experiment, we run the initial exploration phase

of the task proposer with the same prompts used in the main experiments (shown in the last section), and we consider a website to be marked positive for unsafe content if the task proposer generates “N/A” rather than a task. We then calculate accuracy, precision, and recall for various LLMs.

E.2 DETAILS FOR RELIABILITY EXPERIMENTS

We evaluated the verifiable rate of tasks generated by the initial phase of task generation manually. For a set of 100 randomly sampled websites marked as safe by the task proposer in the initial phase, we attempted to complete the generated task, and checked that enough information is present on the website that a solution can be verified. In total, we annotated 300 tasks for Table 2 in 6 hours, and provide the 100 website URLs used in this experiment in the following code block.

```
reliability_sites_list = ['godaddy.com', 'chrome.google.com', '
    → apple.com', 'support.cloudflare.com', 'support.apple.com', '
    → edition.cnn.com', 'go.microsoft.com', 'google.de', 'w3.org',
    → 'yandex.ru', 'bfdi.bund.de', 'microsoft.com', 'apps.apple.
    → com', 'networksolutions.com', 'support.mozilla.org', 'yelp.
    → com', 'cnn.com', 'ec.europa.eu', 'developer.mozilla.org', '
    → icann.org', 'books.google.com', 'globeonewswire.com', '
    → onlinelibrary.wiley.com', 'gnu.org', 'slideshare.net', '
    → metacpan.org', 'porkbun.com', 'oag.ca.gov', 'spiegel.de', '
    → linuxfoundation.org', 'help.opera.com', 'mayoclinic.org', '
    → podcasts.apple.com', 'nhs.uk', 'addons.mozilla.org', 'google.
    → fr', 'pewresearch.org', 'finance.yahoo.com', 'weforum.org',
    → 'g2.com', 'savethechildren.org', 'news.com.au', 'biblia.com
    → ', 'yr.no', 'engadget.com', 'microsoftstore.com', 'ema.
    → europa.eu', 'theintercept.com', 'princeton.edu', '
    → foodandwine.com', 'sfgate.com', 'voguebusiness.com', '
    → ourworldindata.org', 'livingwage.org.uk', 'cms.law', '
    → msdmanuals.com', 'websitesetup.org', 'support.xbox.com', '
    → treehugger.com', 'tripadvisor.com.pe', 'mondragon.edu', '
    → greenparty.ca', 'aaojournal.org', 'restaurantpassion.com', '
    → iwillteachyoutoberich.com', 'moneyconvert.net', '
    → gesundheitsinformation.de', 'ovc.uoguelph.ca', 'zdnet.be', '
    → oxfordamerican.org', 'snackandbakery.com', 'journals.uic.edu
    → ', 'confused.com', 'standards.globalspec.com', '
    → onlyinyourstate.com', 'ahsgardening.org', 'wyze.com', '
    → nornickel.ru', 'viessmann.fr', 'benetton.com', 'firecomm.gov.
    → mb.ca', 'executedtoday.com', 'eukn.eu', 'fraeylemaborg.nl',
    → 'verizon.com/about/news-center', 'orthodoxalbania.org', '
    → cheapjoes.com', 'bake-eat-repeat.com', '
    → plattformpatientensicherheit.at', 'hifinews.com', '
    → cellsignal.com', 'thenotariessociety.org.uk', 'chosenfoods.
    → com', 'westerndressageassociation.org', 'pridesource.com', '
    → northtacomapediatricdental.com', 'strade-bianche.it', '
    → pvdairport.com', 'institute.sandiegozoo.org', 'raintaxi.com
    → ']
```

E.3 AUTOMATIC TASK CATEGORIZATION

We employ *Llama 3.1 70B* to categorize tasks. We prompt *Llama 3.1 70B* with the system prompt in Figure 12 to assign a category in 3 words or less to encourage simple categories. Categories have 16.9 tasks on average, and 953 categories have more than the mean, while 7741 have less than the mean. There is occasional overlap between categories, which can be observed in Figure 13, but for the purposes of understanding performance by category, overlap is acceptable provided categories have sufficiently large numbers of tasks, and performance per category can be accurately calculated. We provide our task categorization script in the official code release.


```

You are a helpful scientific assistant categorizing tasks on the
  → web. You will observe a domain and web navigation task, and
  → you should provide a concise categorization of the task in 3
  → words or less. For example, if the domain is "google.com"
  → and the task is "find a recipe for mashed potato", you may
  → categorize the task as "recipe search".

## Task Format

Here is the format for the task:

[domain]: [task]

Here is what each part means:

`[domain]`: The domain of the website you are observing.
`[task]`: The task a user is trying to accomplish on the website.

## Response Format

Respond with a category name for the task in 3 words or less, and
  → provide only the category name, do not provide an
  → explanation or justification for the categorization.

Here is the next task, please follow the instructions carefully.

```

Figure 12: **System prompt for task categorization.** We employ *Llama 3.1 70B* to automatically label task categories for our dataset. We prompt the LLM to assign categories in 3 words or less, and set the sampling temperature to 0.5 to encourage predictions to use more consistent language. Using these categories, we seek to understand agent performance by category.

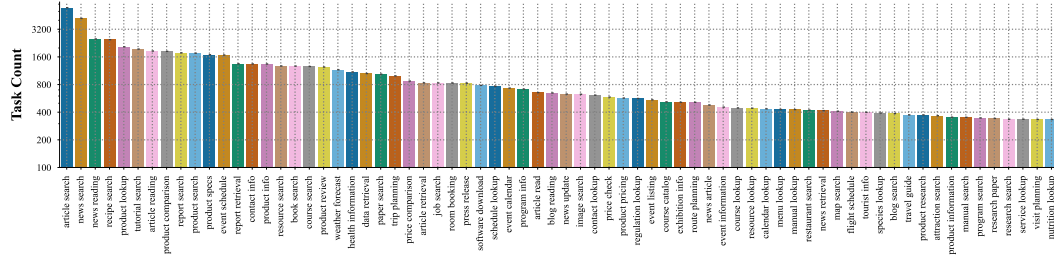


Figure 13: **Largest categories for task generation.** We categorize 150k tasks generated by our pipeline in Section 4, and visualize the number of tasks in the largest 70 categories. Top categories include *article search*, *news search*, *recipe search*, and *product lookup*. The top 12 task categories have more than 1600 tasks assigned to each of them, the mean number of tasks per category is 16.9, and 89% of categories (7741 in total) have fewer than the mean number of tasks.

F UNDERSTANDING AGENT CAPABILITIES & LIMITATIONS

To complement the analyses presented in Section 5, we explore the categories of tasks that agents succeed at most frequently. Shown in Figure 14, we plot the average judge success probability prediction r_T versus task category for the top 70 most successful categories that have at least 100 tasks assigned to them. Based on the figure, top categories include searching for *contact information*, finding *hours of operation*, looking up *biographical information*, obtaining current *weather forecasts*, and conducting *health research*. Based on these results, the top 22 categories are solved with an average success probability > 0.5 using zero-shot agents based on *Llama 3.1 70B*. As stronger

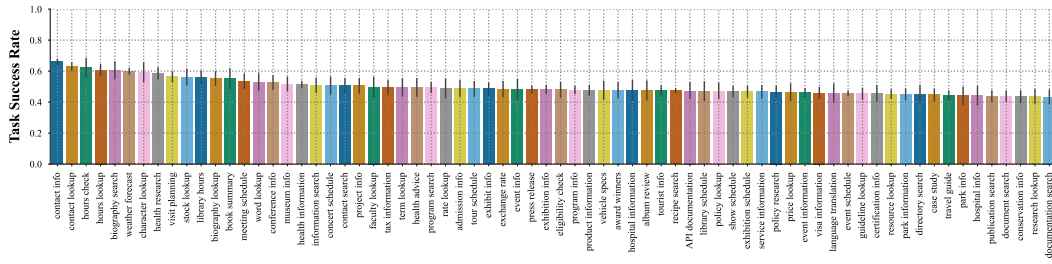


Figure 14: **Most solved categories for task generation.** We explore the completion rates for the top categories of tasks generated by our pipeline. We restrict our focus to categories where at least 100 tasks are assigned, and plots the success rates for the top 70 categories. Results show that 22 of these categories are solved with more than a 50% rate with zero-shot agents based on *Llama 3.1 70B*.

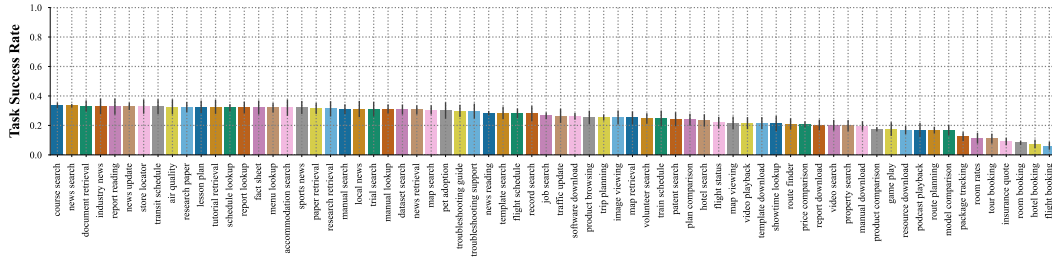


Figure 15: **Least successful categories for internet-scale task generation.** Similar to the previous figure, we explore the rates of task completion for the bottom 70 categories that have at least 100 tasks assigned to them. While the majority of the least successful categories have success rates greater than 20%, performance drops as low as 5%. Many of the categories shown in the plot above involve actions that are not feasible given the current limitations of the Playwright API, and may be possible in future work that extends agents to a fully-operable virtual computer environment. In addition, better LLM backbones are likely to improve performance.

models are developed, the success rates for agents running in our pipeline are likely to improve, and the quality of the data we generate will jointly improve.

In addition to studying the best-performing categories, we also explore the limitations of current agents via their least successful categories. Shown in Figure 15, we select the bottom 70 categories via their average judge success probability for categories with at least 100 tasks assigned. Many of these categories require agents to remember and reason about previous interactions, such as the *product comparison* category. For this category, an agent must review several products, and compare details from memory. In these cases, access to a note-taking tool may improve performance. Additionally, certain task categories involve requests that are not feasible given the limitations of the Playwright API, including categories for *downloading reports / manuals*, and *opening files*. While these tasks are not currently feasible, providing agents with a fully-operable virtual computer environment with applications pre-installed could unlock these abilities in future work.

G AGENT, JUDGE & TASK PROPOSER SYSTEM PROMPTS

We provide the system prompt that powers the agent in this paper. This prompt is released in the code, alongside a fast HTML to Markdown processor we built. The agent prompt is carefully designed to elicit reasoning capabilities, and the experiment in Figure 10 shows the prompt is successful.

```
You are helping me complete tasks by operating a web browser. I
  → will share the current task, and a sequence of webpages and
  → actions from previous steps.

## Action Instructions
```

Based on the information we discovered so far, and the progress we
 ↳ made in previous steps, you are helping me determine the
 ↳ next action.

You will provide an action as JSON in a fenced code block:

```
```json
{
 "action_key": str,
 "action_kwargs": dict,
 "target_element_id": int | null
}
```
```

Actions have the following components:

- 'action_key': The name of the selected action.
- 'action_kwargs': A dictionary of arguments for the action.
- 'target_element_id': An optional id for the element to call the
 ↳ action on.

Action Definitions

I've prepared an API documentation below that defines the actions
 ↳ we can use to complete the task.

Click Action Definition

- 'click': Click on an element specified by 'target_element_id'.

Example Click Action

Suppose you want to click '[id: 5] Sales link':

```
```json
{
 "action_key": "click",
 "action_kwargs": {},
 "target_element_id": 5
}
```
```

Hover Action Definition

- 'hover': Hover over an element specified by 'target_element_id'

Example Hover Action

Suppose you want to hover over '[id: 2] Company Logo image':

```
```json
{
 "action_key": "hover",
 "action_kwargs": {},
 "target_element_id": 2
}
```
```

```

1404 ### Scroll Action Definition
1405
1406 - `scroll`: Scroll the page by `delta_x` pixels to the right and `
1407   ↪ delta_y` pixels down.
1408   - `delta_x`: The number of pixels to scroll to the right.
1409   - `delta_y`: The number of pixels to scroll down.
1410
1411 ### Example Scroll Action
1412
1413 Suppose you want to scroll down the page by 300 pixels:
1414
1415 ```json
1416 {
1417   "action_key": "scroll",
1418   "action_kwargs": {
1419     "delta_x": 0,
1420     "delta_y": 300
1421   },
1422   "target_element_id": null
1423 }
1424 ```
1425
1426 ### Fill Action Definition
1427
1428 - `fill`: Fill an input element specified by `target_element_id`
1429   ↪ with text.
1430   - `value`: The text value to fill into the element.
1431
1432 ### Example Fill Action (Text Input)
1433
1434 Suppose you want to fill `[id: 13] "Name..."` (Enter your name text
1435   ↪ input) with the text `John Doe`:
1436
1437 ```json
1438 {
1439   "action_key": "fill",
1440   "action_kwargs": {
1441     "value": "John Doe"
1442   },
1443   "target_element_id": 13
1444 }
1445 ```
1446
1447 ### Example Fill Action (Range Slider)
1448
1449 Suppose you want to set `[id: 71] "$250 (5)"` (range slider min: 0
1450   ↪ max: 50 step: 1) to the value of `$1000`. The slider has a
1451   ↪ range of 0 to 50 with a step of 1, and the value is
1452   ↪ currently set to `5`. You must translate the desired `$1000`
1453   ↪ to the correct underlying value of `20`:
1454
1455 ```json
1456 {
1457   "action_key": "fill",
1458   "action_kwargs": {
1459     "value": "20"
1460   },
1461   "target_element_id": 71
1462 }

```

```

1458   ``
1459
1460   ### Select Action Definition
1461
1462   - `select`: Select from a dropdown element specified by `
1463     ↪ target_element_id`.
1464     - `label`: The option name to select in the element.
1465
1466   ### Example Select Action
1467
1468   Suppose you want to select the option `red` from `[id: 67] "blue"
1469     ↪ (color select from: red, blue, green)`:
1470
1471   ```json
1472   {
1473     "action_key": "select_option",
1474     "action_kwargs": {
1475       "label": "red"
1476     },
1477     "target_element_id": 67
1478   }
1479   ```
1480
1481   ### Set Checked Action Definition
1482
1483   - `set_checked`: Check or uncheck a checkbox specified by `
1484     ↪ target_element_id`.
1485     - `checked`: Boolean value to check or uncheck the checkbox.
1486
1487   ### Example Set Checked Action
1488
1489   Suppose you want to check `[id: 21] "I agree to the terms and
1490     ↪ conditions" (checkbox)`:
1491
1492   ```json
1493   {
1494     "action_key": "set_checked",
1495     "action_kwargs": {
1496       "checked": true
1497     },
1498     "target_element_id": 21
1499   }
1500   ```
1501
1502   ### Go Back Action Definition
1503
1504   - `go_back`: Go back to the previous page (`target_element_id`
1505     ↪ must be null).
1506
1507   ### Example Go Back Action
1508
1509   Suppose you want to go back to the previous page:
1510
1511   ```json
1512   {
1513     "action_key": "go_back",
1514     "action_kwargs": {},
1515     "target_element_id": null
1516   }

```

```

1512   ```
1513
1514   ### Goto Action Definition
1515
1516   - `goto`: Navigate to a new page (`target_element_id` must be null
1517     ↪ ).
1518     - `url`: The URL of the page to navigate to.
1519
1520   ### Example Goto Action
1521
1522   Suppose you want to open the DuckDuckGo search engine:
1523
1524   ```json
1525   {
1526     "action_key": "goto",
1527     "action_kwargs": {
1528       "url": "https://www.duckduckgo.com"
1529     },
1530     "target_element_id": null
1531   }
1532   ```
1533
1534   ### Stop Action Definition
1535
1536   - `stop`: Stop when the task is complete, and report your progress
1537     ↪ .
1538     - `answer`: Optional answer sent back to me.
1539
1540   ### Example Stop Action
1541
1542   Suppose the task is complete, and you want to stop and report your
1543     ↪ progress:
1544
1545   ```json
1546   {
1547     "action_key": "stop",
1548     "action_kwargs": {
1549       "answer": "The desired task is now complete."
1550     },
1551     "target_element_id": null
1552   }
1553   ```
1554
1555   ## Formatting Your Response
1556
1557   Write a 200 word revised plan based on new information we
1558     ↪ discovered, and progress we made in previous steps. After
1559     ↪ your response, provide the next action as JSON in a fenced
1560     ↪ code block.

```

1558 We also provide the system prompt used by the judge. The system prompt instructs the judge to
1559 predict JSON within a fenced code block that contains a “success” key, an “efficiency” key, and a
1560 “self_correction” key. The success key represents a score from 0 to 1 that estimates the probability
1561 the task is successfully completed. The efficiency key represents a score from 0 to 1 that estimates
1562 the probability the agent has taken the most efficient path to solve the task. The self correction key
1563 represents a score from 0 to 1 that estimates the probability that the agent has demonstrated self
1564 corrective behaviors during its completion of the task. These behaviors include when the agent
1565 backtracks to a more promising state, re-plans when new information is discovered relevant to the
 task, and recognizes its own mistakes. These are generally behaviors we expect from successful

agents, but for this paper we only filter by the success key to select training data for agents.

```

You are helping me evaluate a browser automation script. I will
  ↪ share a task provided to the script, and a sequence of
  ↪ webpages and actions produced by the script.

## The Action Format

The script produces actions as JSON in a fenced code block:

```json
{
 "action_key": str,
 "action_kwargs": dict,
 "target_element_id": int
}
```

Actions have the following components:

- `action_key`: The name of the selected action.
- `action_kwargs`: A dictionary of arguments for the action.
- `target_element_id`: An optional id for the element to call the
  ↪ action on.

## Action Definitions

I've prepared an API documentation below that defines the actions
  ↪ the script can use to complete the task.

### Click Action Definition

- `click`: Click on an element specified by `target_element_id`.

### Example Click Action

Here is an example where the script clicks `[id: 5] Sales link`:

```json
{
 "action_key": "click",
 "action_kwargs": {},
 "target_element_id": 5
}
```

### Hover Action Definition

- `hover`: Hover over an element specified by `target_element_id`

### Example Hover Action

Here is an example where the script hovers over `[id: 2] Company
  ↪ Logo image`:

```json
{
 "action_key": "hover",

```

```

1620 "action_kwargs": {},
1621 "target_element_id": 2
1622 }
1623 ```
1624
1625 ### Scroll Action Definition
1626
1627 - `scroll`: Scroll the page by `delta_x` pixels to the right and `
1628 ↪ delta_y` pixels down.
1629 - `delta_x`: The number of pixels to scroll to the right.
1630 - `delta_y`: The number of pixels to scroll down.
1631
1632 ### Example Scroll Action
1633
1634 Here is an example where the script scrolls down the page by 300
1635 ↪ pixels:
1636
1637 ```json
1638 {
1639 "action_key": "scroll",
1640 "action_kwargs": {
1641 "delta_x": 0,
1642 "delta_y": 300
1643 },
1644 "target_element_id": null
1645 }
1646 ```
1647
1648 ### Fill Action Definition
1649
1650 - `fill`: Fill an input element specified by `target_element_id`
1651 ↪ with text.
1652 - `value`: The text value to fill into the element.
1653
1654 ### Example Fill Action (Text Input)
1655
1656 Here is an example where the script fills `[id: 13] "Name..."` (
1657 ↪ Enter your name text input) with the text `John Doe`:
1658
1659 ```json
1660 {
1661 "action_key": "fill",
1662 "action_kwargs": {
1663 "value": "John Doe"
1664 },
1665 "target_element_id": 13
1666 }
1667 ```
1668
1669 ### Example Fill Action (Range Slider)
1670
1671 Here is an example where the script sets `[id: 71] "$250 (5)"` (
1672 ↪ range slider min: 0 max: 50 step: 1) to the value of `$1000
1673 ↪`. This slider has a range of 0 to 50 with a step of 1, and
1674 ↪ the value is currently set to `5`. The script translates the
1675 ↪ desired `$1000` to the correct underlying value of `20`:
1676
1677 ```json
1678 {

```

```

1674 "action_key": "fill",
1675 "action_kwargs": {
1676 "value": "20"
1677 },
1678 "target_element_id": 71
1679 }
1680 ```
1681
1682 ### Select Action Definition
1683
1684 - `select`: Select from a dropdown element specified by `
1685 ↪ target_element_id`.
1686 - `label`: The option name to select in the element.
1687
1688 ### Example Select Action
1689
1690 Here is an example where the script selects the option `red` from
1691 ↪ `[id: 67] "blue" (color select from: red, blue, green)`:
1692
1693 ```json
1694 {
1695 "action_key": "select_option",
1696 "action_kwargs": {
1697 "label": "red"
1698 },
1699 "target_element_id": 67
1700 }
1701 ```
1702
1703 ### Set Checked Action Definition
1704
1705 - `set_checked`: Check or uncheck a checkbox specified by `
1706 ↪ target_element_id`.
1707 - `checked`: Boolean value to check or uncheck the checkbox.
1708
1709 ### Example Set Checked Action
1710
1711 Here is an example where the script checks `[id: 21] "I agree to
1712 ↪ the terms and conditions" (checkbox)`:
1713
1714 ```json
1715 {
1716 "action_key": "set_checked",
1717 "action_kwargs": {
1718 "checked": true
1719 },
1720 "target_element_id": 21
1721 }
1722 ```
1723
1724 ### Go Back Action Definition
1725
1726 - `go_back`: Go back to the previous page (`target_element_id`
1727 ↪ must be null).
1728
1729 ### Example Go Back Action
1730
1731 Here is an example where the script goes back to the previous page
1732 ↪ :

```

```

1728
1729 ```json
1730 {
1731 "action_key": "go_back",
1732 "action_kwargs": {},
1733 "target_element_id": null
1734 }
1735 ```
1736
1737 ### Goto Action Definition
1738
1739 - `goto`: Navigate to a new page (`target_element_id` must be null
1740 ↪).
1741 - `url`: The URL of the page to navigate to.
1742
1743 ### Example Goto Action
1744
1745 Here is an example where the script opens DuckDuckGo search:
1746
1747 ```json
1748 {
1749 "action_key": "goto",
1750 "action_kwargs": {
1751 "url": "https://www.duckduckgo.com"
1752 },
1753 "target_element_id": null
1754 }
1755 ```
1756
1757 ### Stop Action Definition
1758
1759 - `stop`: Stop when the task is complete, and report the progress.
1760 - `answer`: Optional answer from the script.
1761
1762 ### Example Stop Action
1763
1764 Here is an example where the script stops and reports its progress
1765 ↪ :
1766
1767 ```json
1768 {
1769 "action_key": "stop",
1770 "action_kwargs": {
1771 "answer": "The desired task is now complete."
1772 },
1773 "target_element_id": null
1774 }
1775 ```
1776
1777 ## Evaluation Instructions
1778
1779 Based on the progress of the script, you are helping me determine
1780 ↪ if the desired task has been completed successfully.
1781
1782 You will provide scores as JSON in a fenced code block:
1783
1784 ```json
1785 {
1786 "success": float,

```

```

1782 "efficiency": float,
1783 "self_correction": float
1784 }
1785 ```
1786
1787 ### Score Definitions
1788
1789 - `success`: Your confidence the desired task has been completed
1790 ↳ successfully.
1791 - range: 0.0 (not possible) to 1.0 (absolutely certain).
1792
1793 - `efficiency`: Your confidence the script has taken the most
1794 ↳ efficient path to complete the task.
1795 - range: 0.0 (not possible) to 1.0 (absolutely certain).
1796
1797 - `self_correction`: Your confidence the script has demonstrated
1798 ↳ self-corrective behaviors during its completion of the task.
1799 ↳ These behaviors include backtracking to a more promising
1800 ↳ state, replanning when new information is discovered, and
1801 ↳ recognizing its own mistakes.
1802 - range: 0.0 (not possible) to 1.0 (absolutely certain).
1803
1804 Write a 300 word analysis that establishes specific criteria to
1805 ↳ rigorously evaluate whether the task was completed, followed
1806 ↳ by which criteria the script has satisfied. After your
1807 ↳ response, provide your scores as JSON in a fenced code block.
1808 ↳

```

Finally, we provide the system prompt used in the task proposer to refine the task generated by the first iteration, and to raise the difficulty. The task proposer is instructed via the system prompt to produce JSON with a “proposed\_task” key that represents the task for the agent to complete, a “steps” key that represents the steps that an agent would follow to complete the task, and a “criteria” key that represents the criteria the judge will employ to determine if the task has been completed.

```

1813
1814 You are helping me instruct a language model agent that interacts
1815 ↳ with and navigates live webpages. We instructed the agent to
1816 ↳ complete an initial task, and I will share a sequence of
1817 ↳ webpages visited by the agent during its operation.
1818
1819 ## Your Instructions
1820
1821 Help me refine the task, steps and criteria to raise the
1822 ↳ difficulty, while balancing the agent’s capacity to
1823 ↳ successfully complete the task.
1824
1825 You will provide a task as JSON in a fenced code block:
1826
1827 ```json
1828 {
1829 "proposed_task": str,
1830 "steps": List[str],
1831 "criteria": List[str]
1832 }
1833 ```
1834
1835 Tasks have the following components:
1836
1837 - `proposed_task`: Instruct the agent to complete a task for you
1838 ↳ as if you are a real user that wants help on the website.

```

```

1836 - `steps`: Precise steps in an efficient trajectory that completes
1837 ↪ the task.
1838 - `criteria`: Ground truth answers and criteria to determine if
1839 ↪ the agent completes the task.
1840
1841 Tasks must adhere to the following guidelines:
1842
1843 - Must not require logging in, or making an account.
1844 - Must not require making a purchase, booking, or placing an order
1845 ↪ .
1846 - Must not require creating, deleting, or modifying any posts,
1847 ↪ articles, or webpages.
1848
1849 ## Example Tasks
1850
1851 I've prepared some examples to inspire your task design.
1852
1853 ### `liveevents.iadb.org`
1854
1855 In this example, we explored `liveevents.iadb.org` and saw an
1856 ↪ event page for the IDB Annual Meetings, which includes a
1857 ↪ list of the official hotels and instructions for official
1858 ↪ delegations.
1859
1860 ```json
1861 {
1862 "proposed_task": "I'm attending the IDB Annual Meetings and
1863 ↪ need to find accommodation. Please provide the address
1864 ↪ and phone number for the 'Pullman Santiago Vitacura' and
1865 ↪ 'Double Tree by Hilton' hotels. Additionally, what
1866 ↪ specific details do official delegations need to provide
1867 ↪ to access their special hotel block?",
1868 "steps": [
1869 "Navigate to 'https://liveevents.iadb.org'.",
1870 "Click on the 'Hotels' link in the navigation menu.",
1871 "Locate 'Pullman Santiago Vitacura' in the 'OFFICIAL HOTELS
1872 ↪ FOR THE ANNUAL MEETINGS' list and extract its address
1873 ↪ and telephone number.",
1874 "Locate 'Double Tree by Hilton' in the same list and extract
1875 ↪ its address and telephone number.",
1876 "Read the instructions under 'HOTELS FOR OFFICIAL
1877 ↪ DELEGATIONS' to identify the required information for
1878 ↪ accessing the special hotel block.",
1879 "State the addresses and telephone numbers for both hotels
1880 ↪ and the required information for official delegations
1881 ↪ .",
1882],
1883 "criteria": [
1884 "The agent successfully navigates to the 'Hotels' page on '
1885 ↪ liveevents.iadb.org'.",
1886 "The address for Pullman Santiago Vitacura is stated as '
1887 ↪ Avenida Vitacura 3201 Vitacura, 7630578 Santiago,
1888 ↪ Chile'."
1889 "The telephone number for Pullman Santiago Vitacura is
1890 ↪ stated as '+56 2 2944 7800'."
1891 "The address for Double Tree by Hilton is stated as 'Avenida
1892 ↪ Vitacura 2727, Las Condes Santiago, Chile'."
1893 "The telephone number for Double Tree by Hilton is stated as
1894 ↪ '+56 2 2587 7000'."
1895]
1896 }
1897 ```

```

```

1890 "The agent states that official delegations need to include
1891 ↪ 'the name of your country' and 'the code (included in
1892 ↪ the invitation letters to the Governors)' to access
1893 ↪ the special hotel block."
1894]
1895 }
1896 ```
1897
1898 ### `boldtcastle.com`
1899
1900 In this example, we explored `boldtcastle.com` and saw a page with
1901 ↪ information about visiting Boldt Castle, including
1902 ↪ operating dates, admission prices, and how to get to Heart
1903 ↪ Island.
1904
1905 ```json
1906 {
1907 "proposed_task": "Help me plan a visit to Boldt Castle in 2025
1908 ↪ with one adult and one 6-year-old. Please provide the
1909 ↪ operating dates and hours for the 2025 season, the
1910 ↪ admission cost for just the castle for both of us, how to
1911 ↪ get to Heart Island, and the best phone number for
1912 ↪ general inquiries.",
1913 "steps": [
1914 "Navigate to 'boldtcastle.com'.",
1915 "Click on the 'Visiting' link.",
1916 "Click on the 'Plan Your Visit' link.",
1917 "Identify the 2025 season operating dates and hours for
1918 ↪ Boldt Castle.",
1919 "Find the Boldt Castle-only admission price for an adult
1920 ↪ (13+ years).",
1921 "Find the Boldt Castle-only admission price for a child aged
1922 ↪ 6 (5-12 years).",
1923 "Locate information on how to get to Boldt Castle on Heart
1924 ↪ Island.",
1925 "Find the general inquiry phone number for Boldt Castle.",
1926 "Synthesize all collected information into a concise answer
1927 ↪ ."
1928],
1929 "criteria": [
1930 "State the 2025 season operating dates and hours for Boldt
1931 ↪ Castle as May 10 - October 13, 10:30 AM - 6:30 PM.",
1932 "State the adult admission price for Boldt Castle only as
1933 ↪ $13.50.",
1934 "State the admission price for a child aged 6 for Boldt
1935 ↪ Castle only as $9.50.",
1936 "Provide the physical location of Boldt Castle (Heart Island
1937 ↪ , Alexandria Bay, New York) and mention it's only
1938 ↪ accessible by water.",
1939 "Provide the general inquiry phone number for Boldt Castle
1940 ↪ as 315-482-9724."
1941]
1942 }
1943 ```
1944
1945 ### `visitwestchesterny.com`
1946
1947 In this example, we explored `visitwestchesterny.com` and saw a
1948 ↪ page that lists various coffee houses in Westchester County,

```



```

1944 ↪ including their names, addresses, and links to their
1945 ↪ descriptions.
1946
1947 ```json
1948 {
1949 "proposed_task": "Find a cozy coffee shop in Westchester County
1950 ↪ . Navigate to the 'Coffee Houses' section on the Visit
1951 ↪ Westchester NY website. Find a coffee shop described as '
1952 ↪ cozy' and provide its name, full address, and the exact
1953 ↪ sentence from its description that indicates it is cozy
1954 ↪ .",
1955 "steps": [
1956 "Navigate to 'visitwestchesterny.com'",
1957 "Click on 'Things to Do'",
1958 "Click on 'Food and Drink'",
1959 "Click on 'Coffee Houses'",
1960 "Scroll down to view the coffee shop listings.",
1961 "Identify 'Altamira Cafe Bar' (or any other coffee shop)
1962 ↪ described as 'cozy'."],
1963 "Extract the name and address of the identified coffee shop
1964 ↪ .",
1965 "Click the 'Details' link for the identified coffee shop.",
1966 "Identify and extract the exact sentence from the
1967 ↪ description on its dedicated page that indicates it is
1968 ↪ cozy."
1969],
1970 "criteria": [
1971 "Successfully navigate to the 'Coffee Houses' page.",
1972 "Identify a coffee shop described as 'cozy' (e.g., 'Altamira
1973 ↪ Cafe Bar').",
1974 "State the name of the identified coffee shop (e.g., '
1975 ↪ Altamira Cafe Bar').",
1976 "State the full address of the identified coffee shop (e.g.,
1977 ↪ '245 Main St., New Rochelle, NY 10801').",
1978 "Successfully navigate to the 'Details' page for the
1979 ↪ identified coffee shop.",
1980 "Correctly state the exact sentence from the description
1981 ↪ that indicates it is cozy (e.g., 'Relax in the cozy
1982 ↪ shop or take a treat to go with piping hot espresso, a
1983 ↪ cold coffee, delicious desserts and delightful
1984 ↪ sandwiches.')."]
1985 }
1986 ```
1987
1988 ### `odetterestaurant.com`
1989
1990 In this example, we explored `odetterestaurant.com` and saw a '
1991 ↪ Reservations' page, which lists policies for dietary
1992 ↪ accommodations, birthdays, a deposit requirement,
1993 ↪ cancellations, and rescheduling.
1994
1995 ```json
1996 {
1997 "proposed_task": "I want to make a dinner reservation for 4
1998 ↪ people at Odette, and one of my guests has a severe dairy
1999 ↪ allergy. I also want to request a birthday cake for the
2000 ↪ table. What are the key policies I need to be aware of

```

```

 ↪ regarding my guest's allergy, the cake request, and any
 ↪ deposit or cancellation rules for this reservation?",
"steps": [
 "Navigate to 'odetterestaurant.com'",
 "Go to the 'Reservations' page",
 "Identify the policy regarding dairy allergies and other
 ↪ dietary accommodations",
 "Find the policy for requesting a birthday cake, including
 ↪ notice period and cost",
 "Locate the deposit requirement per person for dinner
 ↪ reservations",
 "Determine the cancellation or rescheduling policy and
 ↪ associated timeframe",
 "Synthesize all relevant policies into a concise answer"
],
"criteria": [
 "State that Odette is unable to accommodate guests with
 ↪ dairy allergies or intolerance.",
 "State that cakes require a 72-hour notice and cost SGD78
 ↪ ++.",
 "Confirm a deposit of SGD200 per person is required for
 ↪ dinner reservations.",
 "State that all reservations are final and non-refundable,
 ↪ but changes can be made at least 72 hours prior to the
 ↪ reservation date."
]
}
```


### `dottyabouticecream.co.uk`



In this example, we explored `dottyabouticecream.co.uk` and saw a



- ↪ form for hiring Dotty's ice cream van for corporate events,
- ↪ which includes fields for event details, guest count, and
- ↪ flavor inquiries.



```

```json
{
  "proposed_task": "Inquire about hiring Dotty's ice cream van
    ↪ for a corporate event in Manchester, M1 1AE, on August 15
    ↪ th, 2024, from 2 PM to 4 PM, for 100 guests. Ask if
    ↪ vanilla, honeycomb crunch, and mango sorbet are available.
    ↪ Fill out the 'Get in Touch' form with your details (Jane
    ↪ Doe, jane.doe@example.com, 07123456789) and note you
    ↪ found them via a web search. Do not submit the form.",
  "steps": [
    "Navigate to the 'Get in Touch' page on dottyabouticecream.
      ↪ co.uk.",
    "Fill 'Jane Doe' into the 'Name' field.",
    "Fill 'jane.doe@example.com' into the 'Email' field.",
    "Fill '07123456789' into the 'Telephone Number' field.",
    "Fill 'August 15th, 2024' into the 'Event Date' field.",
    "Fill '2 PM - 4 PM' into the 'Time of Ice Cream Service'
      ↪ field.",
    "Fill 'Manchester, M1 1AE' into the 'Venue Address (incl.
      ↪ Postcode)' field.",
    "Fill '100' into the 'Number of Expected Guests' field.",
    "Fill 'Web Search' into the 'Where Did You Hear About Dotty
      ↪ ?' field.",

```


```

```

2052 "Fill the 'Message' field with an inquiry about the
2053 ↪ availability of 'Vanilla, Honeycomb Crunch, and Mango
2054 ↪ Sorbet' flavors for a corporate event.",
2055 "Confirm all specified fields are accurately filled, but do
2056 ↪ not click the 'Gimmie Ice Cream' submit button."
2057],
2058 "criteria": [
2059 "The agent successfully navigates to the 'Get in Touch' page
2060 ↪ .",
2061 "The 'Name' field is filled with 'Jane Doe'.",
2062 "The 'Email' field is filled with 'jane.doe@example.com'.",
2063 "The 'Telephone Number' field is filled with
2064 ↪ '07123456789'.",
2065 "The 'Event Date' field is filled with 'August 15th,
2066 ↪ 2024'.",
2067 "The 'Time of Ice Cream Service' field is filled with '2 PM
2068 ↪ - 4 PM'.",
2069 "The 'Venue Address (incl. Postcode)' field is filled with '
2070 ↪ Manchester, M1 1AE'.",
2071 "The 'Number of Expected Guests' field is filled with
2072 ↪ '100'.",
2073 "The 'Where Did You Hear About Dotty?' field is filled with
2074 ↪ 'Web Search'.",
2075 "The 'Message' field clearly inquires about the availability
2076 ↪ of 'Vanilla, Honeycomb Crunch, and Mango Sorbet'
2077 ↪ flavors for a corporate event.",
2078 "The agent does not submit the form by clicking the 'Gimmie
2079 ↪ Ice Cream' button."
2080]
2081 }
2082 ""
2083 ### `engineered.polestar.com`
2084 In this example, we explored `engineered.polestar.com` and saw a
2085 ↪ page with information about Polestar Engineered Optimization
2086 ↪ for various Volvo models, including the 2023 Volvo XC60
2087 ↪ with a B5 Drive-E engine.
2088 ```json
2089 {
2090 "proposed_task": "Is a Polestar Engineered Optimization
2091 ↪ available for a 2023 Volvo XC60 with a B5 Drive-E engine?
2092 ↪ If so, what are the primary performance benefits, how
2093 ↪ long does the installation take, and how would I find a
2094 ↪ dealer for this service?",
2095 "steps": [
2096 "Navigate to engineered.polestar.com.",
2097 "Under 'Can My Volvo Be Optimised?', select 'XC' then 'New
2098 ↪ XC60' for the model.",
2099 "Locate and select 'XC60 B5 Drive-E AWD Automatic 2023' or '
2100 ↪ XC60 B5 Drive-E FWD Automatic 2023' to view its
2101 ↪ optimization details.",
2102 "Confirm if the vehicle is 'Approved for Polestar Engineered
2103 ↪ Optimization'.",
2104 "Identify the primary performance benefits listed for the
2105 ↪ optimization.",
2106 "Determine the approximate installation time.",

```

```

2106 "Click on any 'Find a retailer' or 'Contact a dealer' links
2107 ↪ to see where they lead.",
2108 "Based on the website's information, describe how a user
2109 ↪ would find a dealer for installation.",
2110 "Synthesize all gathered information to answer the task."
2111],
2112 "criteria": [
2113 "Confirm that a 2023 Volvo XC60 with a B5 Drive-E engine is
2114 ↪ 'Approved for Polestar Engineered Optimization'.",
2115 "State the primary performance benefit as 'Power Mid-Range
2116 ↪ up to (hp) +3%' (from the specific product page) or '
2117 ↪ Up to +15% increased mid-range power' (from the
2118 ↪ general 'Get optimisation' page).",
2119 "State that the installation takes 'less than 60 minutes'.",
2120 "Clearly state that clicking the 'Find a retailer' or '
2121 ↪ Contact a dealer' links does not lead to a functional
2122 ↪ dealer search tool, and that users are advised to
2123 ↪ contact their local Volvo retailer directly for
2124 ↪ further questions."
2125]
2126 }
2127 \ \ \
2128 ### `ajga.org`
2129
2130 In this example, we explored `ajga.org` and saw a page with
2131 ↪ information about Performance Based Entry (PBE) Stars for
2132 ↪ junior golfers, including how they carry over to the next
2133 ↪ season and tips for maximizing tournament opportunities.
2134
2135 ```json
2136 {
2137 "proposed_task": "As a junior golfer, help me understand how my
2138 ↪ Performance Based Entry (PBE) Stars carry over to the
2139 ↪ next season on the AJGA circuit. Additionally, I'm
2140 ↪ looking for two tips from the website on how to best
2141 ↪ maximize my opportunities to play in tournaments. Can you
2142 ↪ provide this information?",
2143 "steps": [
2144 "Navigate to the 'Juniors' section of the AJGA website.",
2145 "Click on 'How to Play in the AJGA'.",
2146 "Click on the 'PBE' link within the 'How to Play in the AJGA
2147 ↪ ' section.",
2148 "Scroll down to locate and click on the 'Important Notes and
2149 ↪ Tips' link.",
2150 "Identify and extract information regarding the carry-over
2151 ↪ of Performance Stars to the next season.",
2152 "Identify and extract two distinct tips that help players
2153 ↪ maximize their tournament playing opportunities.",
2154 "Synthesize the gathered information to provide a
2155 ↪ comprehensive answer."
2156],
2157 "criteria": [
2158 "Explain that Performance Stars earned in one year (e.g.,
2159 ↪ 2024) are carried over to the beginning of the next
2160 ↪ season (e.g., 2025) and are combined with new
2161 ↪ membership Performance Stars based on grad year.",

```

```

2160 "Identify and state that 'Plan your tournament schedule
2161 ↪ early to maximize playing opportunities and prevent
2162 ↪ missed deadlines.' is a key recommendation.",
2163 "Identify and state that 'Qualifiers are great opportunities
2164 ↪ for all players to earn Performance Stars and build
2165 ↪ their status.' is a second key recommendation."
2166]
2167 }
2168 ```
2169 ### `passports.gov.au`
2170
2171 In this example, we explored `passports.gov.au` and saw a page
2172 ↪ with a section about the documents needed to prove
2173 ↪ Australian citizenship for individuals born in Australia on
2174 ↪ or after August 20, 1986.
2175
2176 ```json
2177 {
2178 "proposed_task": "I was born in Australia on or after August
2179 ↪ 20, 1986, and am applying for my first Australian
2180 ↪ passport. What documents do I need to prove my Australian
2181 ↪ citizenship? Please list the primary document,
2182 ↪ acceptable alternatives, and any specific requirements
2183 ↪ for proving citizenship by birth based on my parents' or
2184 ↪ grandparents' status.",
2185 "steps": [
2186 "Navigate to 'passports.gov.au'.",
2187 "Navigate to the 'How it works' section.",
2188 "From 'How it works', navigate to 'Documents you need'.",
2189 "On the 'Documents you need' page, navigate to the '
2190 ↪ Citizenship' section.",
2191 "Within the 'Citizenship' section, locate the information
2192 ↪ for individuals 'Born in Australia on or after 20
2193 ↪ August 1986'."
2194],
2195 "criteria": [
2196 "Identify the primary document required for proof of
2197 ↪ citizenship.",
2198 "Identify and list all acceptable alternative documents.",
2199 "Detail the specific scenarios for proving citizenship by
2200 ↪ birth, including those involving parents' or
2201 ↪ grandparents' documentation and the special case for
2202 ↪ permanent resident parents."
2203],
2204 "criteria": [
2205 "State that the primary document is the applicant's full
2206 ↪ Australian birth certificate.",
2207 "List an Australian citizenship certificate in the applicant
2208 ↪ 's name as an acceptable alternative.",
2209 "List an Australian passport issued in the applicant's name
2210 ↪ on or after 1 January 2000 that was valid for at least
2211 ↪ two years as an acceptable alternative.",
2212 "Detail the scenario where one parent was an Australian
2213 ↪ permanent resident or citizen, specifying the required
2214 ↪ parental documents (birth certificate, passport, or
2215 ↪ citizenship certificate).",
2216 "Explicitly mention that if both parents were Australian
2217 ↪ permanent residents when the applicant was born,
2218 ↪ evidence of citizenship must be obtained from the
2219 ↪ Department of Home Affairs."
2220]
2221 }

```

```

 "Include the scenario involving grandparents' documents (
 ↳ birth certificate, passport, or citizenship
 ↳ certificate) if the parent was born in Australia on or
 ↳ after 20 August 1986."
]
}
'''

Formatting Your Response

Establish how the task can be refined in at most 300 words, and
 ↳ synthesize relevant content and features on the website in
 ↳ your response. After your response, provide a refined task
 ↳ as JSON in a fenced code block.

```

## H DETAILS FOR TRAINING AGENTS

To understand the utility of data we obtained, we train agents and test on four relevant benchmarks: InSTA, WebVoyager (He et al., 2024), Mind2Web (Deng et al., 2023), WebLINX (Lù et al., 2024). In particular, our test set consists of a held-out set of 3,000 websites and tasks produced by the task generation feedback loop. Note these websites are not present in the training set. For WebVoyager (He et al., 2024), we transfer agents trained on our data zero-shot to 643 tasks on 15 websites WebVoyager (He et al., 2024). The websites in the WebVoyager benchmark are not present in the 20k trajectories we collected in Section 6.1 we used for training agents. For this experiment, we fine-tuned models based on *Qwen 3 1.7B* with a maximum sequence length of 16,384 tokens, and the most recent 5 observations, and actions in the context. We employed full fine-tuning on this model, with Adam, learning rate of  $5e-5$ , batch size of 32, `bfloat16`, and other parameters kept as the PyTorch defaults for Adam. Each model was trained using one epoch, a linear warm-up corresponding to the first 0.01 steps of training, and a linear decay to  $6e-5$  afterwards.

To filter data, we select trajectories that were scores as `Judge(Success) = 1`, which corresponded to 10.5k of 20k trajectories produced by the *Qwen 3 235B* data collection policy. Scores for filtering were produced by a *Qwen 3 235B* judge. We employed a simple filtering strategy that only considers the success score from the judge, and no other filtering conditions were used. Note the judge also predicts efficiency, and self correction scores, which could likely also help select the best data for training, but we did not explore filtering by these scores in this work.

For experiments on static benchmarks, we fine-tune `google/flan-t5-large` for Mind2Web, and `meta-llama/Llama-3.1-8B-Instruct` for WebLINX using official fine-tuning code released with corresponding benchmarks. We employ identical training hyperparameters to those used by Lù et al. (2024) for Llama in their official training code and Deng et al. (2023) for Flan to ensure that our results are comparable to previous work. Section 6.2 reports performance on the official `test_web` split of the WebLINX benchmark, and the official `test_website` split of the Mind2Web benchmark, where agents are tested on unobserved websites. The websites in these static benchmarks were not present in the dataset we generated for this experiment to ensure fairness.

## I HYPERPARAMETERS

We provide a list of the hyperparameters used in this work in Table 3. Hyperparameters for mixing our data with human data on static benchmarks are selected to mirror prior work in synthetic data (Trabucco et al., 2024), and to adhere to standard hyperparameters for WebLINX (Lù et al., 2024), and Mind2Web (Deng et al., 2023). We train using all available human data on these benchmarks, and add trajectories filtered using the previously discussed methodology, sampled at a 20% rate in the data-loader compared to an 80% rate for human data.

| Hyperparameter Name                               | Value                                                                                                                                                                                                                                       |
|---------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Models Used For Agents                            | Qwen/Qwen3-1.7B<br>Qwen/Qwen3-235B-A22B<br>meta-llama/Llama-4-Maverick-17B-128E-Instruct<br>meta-llama/Llama-3.1-70B-Instruct<br>meta-llama/Llama-3.3-70B-Instruct<br>google/gemini-2.5-flash                                               |
| Models Used For Judges                            | Qwen/Qwen3-235B-A22B<br>meta-llama/Llama-4-Maverick-17B-128E-Instruct<br>meta-llama/Llama-3.1-70B-Instruct<br>meta-llama/Llama-3.3-70B-Instruct<br>google/gemini-2.5-flash<br>google/gemini-1.5-pro<br>openai/gpt-4.1-mini<br>openai/gpt-4o |
| Common Crawl PageRank                             | cc-main-2024-apr-may-jun-host-ranks.txt.gz                                                                                                                                                                                                  |
| Number of sites before filtering                  | 1,000,000                                                                                                                                                                                                                                   |
| Number of tasks after filtering                   | 146,746                                                                                                                                                                                                                                     |
| Max Tokens Per Observation                        | 2,048                                                                                                                                                                                                                                       |
| Max Tokens Per Agent Trace                        | 1,024                                                                                                                                                                                                                                       |
| Max Tokens Per Judge Trace                        | 1,024                                                                                                                                                                                                                                       |
| Max Tokens Per Task Proposer Trace                | 1,024                                                                                                                                                                                                                                       |
| Last Steps Per Agent Context                      | 5                                                                                                                                                                                                                                           |
| Last Steps Per Judge Context                      | 5                                                                                                                                                                                                                                           |
| Last Steps Per Task Proposer Context              | 5                                                                                                                                                                                                                                           |
| Task Proposer Feedback Loops                      | 1                                                                                                                                                                                                                                           |
| LLM Sampling Temperature                          | 0.5                                                                                                                                                                                                                                         |
| LLM Sampling Top P                                | 1.0                                                                                                                                                                                                                                         |
| LLM Sampling Top K                                | default                                                                                                                                                                                                                                     |
| Fine-tuned LLM in Section 6.1                     | Qwen/Qwen3-1.7B                                                                                                                                                                                                                             |
| InSTA Training Epochs                             | 1                                                                                                                                                                                                                                           |
| InSTA Batch Size                                  | 32                                                                                                                                                                                                                                          |
| InSTA Learning Rate                               | 5e-5                                                                                                                                                                                                                                        |
| InSTA Optimizer                                   | Adam                                                                                                                                                                                                                                        |
| Mind2Web LLM                                      | google/flan-t5-large                                                                                                                                                                                                                        |
| Mind2Web Training Iterations                      | 11,505                                                                                                                                                                                                                                      |
| Mind2Web Batch Size                               | 32                                                                                                                                                                                                                                          |
| Mind2Web Learning Rate                            | 5e-5                                                                                                                                                                                                                                        |
| Mind2Web Optimizer                                | Adam                                                                                                                                                                                                                                        |
| WebLINX LLM                                       | meta-llama/Llama-3.1-8B-Instruct                                                                                                                                                                                                            |
| WebLINX Training Iterations                       | 10,000                                                                                                                                                                                                                                      |
| WebLINX Batch Size                                | 16                                                                                                                                                                                                                                          |
| WebLINX Learning Rate                             | 5e-5                                                                                                                                                                                                                                        |
| WebLINX Optimizer                                 | Adam                                                                                                                                                                                                                                        |
| Data Filtering Condition                          | Judge (Success) = 1                                                                                                                                                                                                                         |
| Human Data Sampling Probability $p_{\text{real}}$ | 80%                                                                                                                                                                                                                                         |

**Table 3: Hyperparameters used in our paper.** We organize hyperparameters into seven sections, for the names of LLMs used as agents in the paper, the names of LLMs used as judges in the paper, the hyperparameters used for data collection, the sampling parameters for LLMs, the training parameters for static benchmarks, and the filtering and data mixing hyperparameters.

## J COST ANALYSIS FOR LLAMA 3.1 70B

To understand the significant reduction in cost that we obtain by running LLMs locally to generate data, we analyze the number of tokens processed by the LLM, and compute an expected cost if this were served using proprietary models. As the analysis shows, using *Llama 3.1 70B* is a feasible



option for running agents at this large scale, and results in the paper show that this choice of LLM backbone does not compromise performance. We have deep gratitude to the Llama team at Meta, and the Qwen team at Alibaba for working to make developments in language modeling available to the research community at no cost. We see up to a 95% reduction in cost with these models.

| Variable Name                                                                                             | Value            |
|-----------------------------------------------------------------------------------------------------------|------------------|
| Number of tasks                                                                                           | 146,746          |
| Average tokens per observation                                                                            | 1,024            |
| Max observations per agent context window                                                                 | 5                |
| Average agent / judge response size                                                                       | 512              |
| Max tokens per system prompt                                                                              | 1,024            |
| Average steps per task                                                                                    | 15               |
| Estimated tokens processed by the agent                                                                   | 14.65B tokens    |
| Tokens processed by the judge                                                                             | 1.35B tokens     |
| Total tokens processed                                                                                    | 16.00B tokens    |
| Expected API cost for <i>GPT-4.1</i>                                                                      | \$ 32,000.00     |
| Expected API cost for <i>Gemini 2.5 Pro</i>                                                               | \$ 20,000.00     |
| Expected AWS compute cost for serving <i>Llama 3.1 70B</i><br>(3,840 v100 GPU hours using spot instances) | \$ 1,575.70      |
| Percent saved using <i>Llama 3.1 70B</i>                                                                  | [95.08, 92.12] % |

Table 4: **Cost analysis for different LLM models in the fully-scaled pipeline.** This table provides statistics for the number of tokens that were processed by our pipeline, and why serving using a local LLM engine like vLLM is important for bringing down costs.