

AGENT Q: ADVANCED REASONING AND LEARNING FOR AUTONOMOUS AI AGENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) have shown remarkable capabilities in natural language tasks requiring complex reasoning, yet their application in agentic, multi-step reasoning within interactive environments remains a difficult challenge. Traditional supervised pre-training on static datasets falls short in enabling autonomous agent capabilities needed to perform complex decision-making in dynamic settings like web navigation. Previous attempts to bridge this gap through supervised fine-tuning on curated expert demonstrations often suffer from compounding errors and limited exploration data, resulting in sub-optimal policy outcomes. To overcome these challenges, we propose a framework that combines guided Monte Carlo Tree Search (MCTS) search with a self-critique mechanism and iterative fine-tuning on agent interactions using an off-policy variant of the Direct Preference Optimization (DPO) algorithm. Our method allows LLM agents to learn effectively from both successful and unsuccessful trajectories, thereby improving their generalization in complex, multi-step reasoning tasks. We validate our approach in the WebShop environment, a simulated e-commerce platform—where it consistently outperforms behavior cloning and reinforced fine-tuning baseline, and **beats average human performance** when equipped with the capability to do online search. In real-world booking scenarios, our methodology boosts Llama-3 70B model’s zero-shot performance from **18.6% to 81.7%** success rate (a **340% relative increase**) after a single day of data collection and further to **95.4%** with online search. We believe this represents a substantial leap forward in the capabilities of autonomous agents, paving the way for more sophisticated and reliable decision-making in real-world settings.

1 INTRODUCTION

The recent advances in Large Language Models (LLMs) represent a significant leap in artificial intelligence. Frontier models like ChatGPT (John Schulman et al., 2022), Gemini (Anil et al., 2023), Opus (Anthropic, 2024), and LLaMA-3 (Touvron et al., 2023) demonstrate promising reasoning capabilities that approach average human performance in a number of domains. These breakthroughs have extended the utility of LLMs from traditional chat and text-based applications to more dynamic, agentic roles, in which they do not just generate text but can take actions autonomously in a number of environments including code and software engineering (Holt et al., 2024; Zhang et al., 2024d; Jimenez et al., 2024; Yang et al., 2024), device control (Wang et al., 2024a; Zhang et al., 2023; Chen and Li, 2024) and web applications (Hong et al., 2023; Deng et al., 2023; Zhou et al., 2024b; Lai et al., 2024a; Gur et al., 2024) among others. However, despite these advancements, significant challenges persist: LLMs still struggle to generalize effectively in interactive, multi-step environments, since they are not native trained for such applications. This is true, even for some of the strongest models of the current generation, such as GPT-4 (Achiam et al., 2023).

A growing literature on agentic formulation seeks to address these issues; however these works mostly focus on building frameworks around prompt-based learning on existing models or limited fine-tuning on static datasets, and are thus limited by the base models’ reasoning and decision making capabilities. Reasoning and planning have indeed been highlighted as core challenges for current LLMs. Since the seminal work on chain-of-thought reasoning (Wei et al., 2022), significant efforts have been made to improve these capabilities via prompt-based strategies (Kojima et al., 2022; Wang et al., 2023; Qiao et al., 2023; Yao et al., 2023a). While successful, these approaches are still bounded

054 by the base model’s performance. Another direction of research has explored fine-tuning approaches
055 (Zelikman et al., 2022; Pang et al., 2024), and more recently combining them with inference-time
056 search prompting (Yao et al., 2023a) to produce fine-grained feedback. Concurrent works (Xie et al.,
057 2024; Hwang et al., 2024; Zhang et al., 2024e; Tian et al., 2024) utilize the traces produced by search
058 algorithms and combine them with optimization approaches (Rafailov et al., 2023; Zelikman et al.,
059 2022) to achieve significant boost in capabilities, especially in mathematics problem solving and
060 code generation.

061 In this work we explore improving planning and reasoning capabilities of a web agent, which interacts
062 with a real world website. Our goal is to design an approach that allows the agent to improve with
063 autonomous experience and limited supervision. Indeed, prior works (Yao et al., 2023b; Zhang et al.,
064 2024c; Masterman et al., 2024; Sumers et al., 2024) have shown strong reasoning to be critical for
065 performance of autonomous agents, where challenges are even greater than during text generation,
066 as the model needs to further understand how its actions affect its environment. Towards this goal,
067 we introduce **Agent Q**—a novel approach that combines several key concepts in reasoning, search,
068 self-critique and reinforcement learning. Our method takes inspiration from Sutton’s The Bitter
069 Lesson on the power of general purpose methods that continue to scale with increased computation,
070 showing the significant benefits of combining *search* and *learning*.

071 Inspired by the success of search-based methods in prior game-playing settings (Silver et al., 2017a;
072 Brown and Sandholm, 2019; Gray et al., 2021) and mathematical reasoning (Yao et al., 2023a; Besta
073 et al., 2024), we deploy a Monte Carlo Tree Search (MCTS) based search routine over web pages to
074 guide agent exploration. Given the complexity of the environment, we use a base LLM for sampling
075 possible rationales and web actions to explore. While this simple search-strategy shows a meaningful
076 improvement in the success rate, it still struggles on long horizon tasks due to sparsity of environment
077 rewards. Indeed even a small mistake across the trajectory can cause the final agent output to be
078 wrong, creating significant credit assignment problems. To overcome this, we use AI feedback (Bai
079 et al., 2022) and self-criticism (Yuan et al., 2024) to further prompt the LLM to provide self-evaluation
080 feedback at each node, which serves as intermediate reward and helps guide the search steps. This
081 meaningfully improves the final agent success rate, but requires significant online interactions and
082 moreover the capability to rollback actions, which is not always possible in online realistic settings.
083 Such online autonomous search with little supervision on the web can result in a weak or unsafe
084 agent which can make many errors, resulting in risky behaviors in sensitive online settings like bank
transfers and sensitive information sharing.

085 To correct this, we use the traces generated by the search process to improve capabilities of the
086 model by learning from both the successful and unsuccessful trajectories with offline reinforcement
087 learning, utilizing the Direct Preference Optimization (DPO) algorithm. We create preferences over
088 different branches at the node level, which are scored using a mixture of the AI process feedback
089 rewards and the final success rate of the explored branch. We evaluate our approach on the simulated
090 WebShop benchmark (Yao et al., 2022)—a simulated e-commerce platform—as well as a real-world
091 reservations booking website. We utilize LLaMa 3-70B as the base model in our experiments. In the
092 WebShop environment, our approach consistently outperforms behavior cloning and reinforcement
093 learning fine-tuned baselines, and **beats average human performance when equipped with the
094 capability to do online search.**

095 In our real-world booking experiments, using our **Agent Q** framework we improve the model zero-
096 shot absolute success rate from **18.6% to 81.7%** (a **340% relative increase**), outperforming GPT-4’s
097 performance after a single day of autonomous data collection. When we equip Agent Q with online
098 search capability, our absolute success further improves to **95.4%**. We believe that our approach
099 represents a significant step forward in the development of autonomous web agents through it’s search
100 and self-critique capabilities, setting a new benchmark for reliable multi-step decision-making in
101 interactive settings.

102 2 RELATED WORK

103 Our work touches on a large number of research directions around agent design, self-improvement,
104 reasoning and reinforcement learning. We include a short overview of related works from those
105 various fields below.

2.1 GUIDED SEARCH FOR REASONING AND PLANNING

The latest generation of Large Language Models (LLMs) have demonstrated promising emerging properties around reasoning and planning. Moreover such behaviours can be directly elicited from strong models only using simple prompting techniques (Wei et al., 2022; Kojima et al., 2022; Qiao et al., 2023). These have also become an integral part of agentic design (Yao et al., 2023b; Zhang et al., 2024c), which we also utilize for our approach. Another emerging research direction is based around step-by-step verifiers or “Process Reward Models” (Uesato et al., 2022; Lightman et al., 2023), specifically for mathematical reasoning. These have shown to improve performance beyond purely outcome-based training, however they require a large amount of human effort to label individual steps. Some recent approaches have proposed self-supervised methods for step-level supervision (Hwang et al., 2024; Wang et al., 2024b; Setlur et al., 2024a). A number of concurrent works (Xie et al., 2024; Zhang et al., 2024e; Tian et al., 2024) have further explored tree-based search approaches (Yao et al., 2023a) in combination with DPO (Rafailov et al., 2023) training for math-based reasoning. These algorithms optimize actions at the node level, using different branches produced by the search algorithm to create preference pairs. Our approach shares similarities to the self-supervised search proposed in (Yao et al., 2023a) with a combination of AI-based feedback (Bai et al., 2022; Yuan et al., 2024) to guide intermediate search steps, but we are the first to scale this a realistic agent setting. Similar approaches were proposed in (Zhou et al., 2024a; Hao et al., 2023; Kang et al., 2024), and other works (Koh et al., 2024); however these works only use the base model’s zero-shot capability to search and do not train it further. Moreover they are only evaluated on simulated environments. Beyond the search stage, our work further adopts the training methodology of (Xie et al., 2024; Zhang et al., 2024e; Tian et al., 2024), which significantly boosts our agent’s zero-shot capabilities.

2.2 WEB AGENTS

The strength and capabilities of recent pretrained Large Language (Vision) Models LL(V)Ms has significantly boosted progress in developing autonomous web-agents. Improved code understanding and long context have allowed agents to represent environment state and action space with document object model (DOM) allowing for deployment in complex and realistic domains. Moreover strong reasoning (Yao et al., 2023b) and planning (Liu et al., 2023; Zhang et al., 2024c) capabilities have also led to the development of a number of promising agents (Zhang and Zhang, 2023; Hong et al., 2023; Zhou et al., 2024b; Deng et al., 2023; Gur et al., 2024). Beyond using LL(V)Ms as plug-and-play planners/policies, recent works have sought to improve agentic-specific performance. Examples include online exploration (Zhang et al., 2024a), planning (Zhang et al., 2024b), error-correction (Wang et al., 2024a), and self- (Wu et al., 2024) or AI-critique (He et al., 2024; Pan et al., 2024). However, with small exceptions (Nakano et al., 2022) (which is still limited in scope) these agents mostly provide a framework around a strong pre-existing model like GPT4-V or deploy limited fine-tuning and adaptation. In this work we show that model training is crucial for continuous improvement. We combine a planning and reasoning agent with MCTS inference-time search and AI self-critique for self-supervised data collection, which we then use for RL type training.

2.3 REINFORCEMENT LEARNING FOR LLMs AND AGENTS

Reinforcement Learning has become a significant component of training modern generative AI systems (Ouyang et al., 2022; Bai et al., 2022; Touvron et al., 2023). Classical approaches have deployed the PPO algorithm (Schulman et al., 2017)—or similar policy-gradient based methods—and have even been scaled to autonomous web search agents (Nakano et al., 2022) as well as embodied applications with vision-language models (Zhai et al., 2024) (in simulation). However, these algorithms are challenging due to their complexity and the need for a high number of online samples from the model. This is especially prominent in potentially risky situations, such as autonomous agentic models that could make a number of impactful mistakes during training. Implicit Language Q-learning (Snell et al., 2022) and the Q-transformer (Chebotar et al., 2023) are offline RL algorithms (Levine et al., 2020) designed for auto-regressive transformer models, and hence can be safely trained on pre-collected datasets; however they have not been successfully scaled to modern LLMs. While these methods represent a token-level MDP, (Zhou et al., 2024c) has shown success formulating the RL problem at a step level and these ideas have recently been scaled to a general device-control agent (Bai et al., 2024). However, these algorithms still have high complexity and require auxiliary models, such as value functions, so instead in our approach we opt to use the Direct

162 Preference Optimization (DPO) algorithm (Rafailov et al., 2023) due to its simplicity and natural fit
 163 for the branching nature of tree-search based data.
 164

165 3 PRELIMINARIES

166
 167 In this section we will outline the preliminaries of our agent training process. For a full description of
 168 our agentic system formulation consider Appendix A. Four training purposes at each time step t the
 169 agent receives a state \mathbf{h}_t and will produce actions $\mathbf{a}_t \sim \pi_{\theta}(\mathbf{a}|\mathbf{h}_t)$.
 170

171 3.1 FINE-TUNING LANGUAGE MODELS FROM FEEDBACK

172
 173 Classical approaches to RLHF in foundation models Stiennon et al. (2022); Ouyang et al. (2022) use
 174 the model as a policy π_{θ} and optimize an objective of the form:

$$175 \mathbb{E}_{\mathbf{a} \sim \pi_{\theta}(\mathbf{a}|\mathbf{h})} [r(\mathbf{a}, \mathbf{h})] - \beta \mathbb{D}_{KL}[\pi_{\theta}(\mathbf{a}|\mathbf{h}) || \pi_{\text{ref}}(\mathbf{a}|\mathbf{h})] \quad (1)$$

176
 177 where π_{ref} is some reference policy (usually the initial model). The goal of this formulation is to
 178 optimize some target objective (expressed by the reward $r(\mathbf{a}, \mathbf{h})$) while preventing out-of-distribution
 179 drift. This objective can be extended to multi-step agentic problems, where the model interacts with
 180 an external environment env such as in Nakano et al. (2021) which focuses on information retrieval
 181 using web navigation. In this case we use an objective of the kind

$$182 \mathbb{E}_{\pi_{\theta}, \text{env}} \left[\sum_t r(\mathbf{a}_t, \mathbf{h}_t) \right] - \beta \mathbb{D}_{KL}[\pi_{\theta}(\mathbf{a}_t|\mathbf{h}_t) || \pi_{\text{ref}}(\mathbf{a}_t|\mathbf{h}_t)] \quad (2)$$

183
 184 Classical RLHF has used policy gradient type of algorithms, such as PPO Schulman et al. (2017),
 185 however, they are complex and require online data, which can be costly/dangerous to collect au-
 186 tonomously in the agent setting. While PPO has shown some success in prior web agent applications
 187 Nakano et al. (2021). The issues above largely make the approach not practical for general web tasks,
 188 beyond information retrieval. In this work we utilize some recent alternatives, outlined below.
 189

190 3.1.1 REINFORCED FINE-TUNING

191
 192 Reinforced fine-tuning (RFT) algorithms Zelikman et al. (2022); Gulcehre et al. (2023); Yuan et al.
 193 (2023); Singh et al. (2024) have grown in popularity due to their simplicity and scalability. These
 194 methods aggregate data and filter out the sub-optimal samples based on some reward model or
 195 a verifier to construct a growing dataset of high-quality trajectories \mathcal{D} . Given this dataset and a
 196 parameterized model π_{θ} we can carry out standard supervised fine-tuning (SFT):

$$197 \mathcal{L}(\pi_{\theta}, \mathcal{D}) = -\mathbb{E}_{\mathcal{D}} \left[\sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t|\mathbf{h}_t) \right] \quad (3)$$

198
 199 In this objective the divergence penalty is only applied implicitly by limiting the number of training
 200 rounds. While simple and relatively successful, empirically these methods tend to under-perform
 201 standard RL and alternatives Dubois et al. (2024); Tajwar et al. (2024); Setlur et al. (2024b) in the
 202 text generation domain, particularly in reasoning. We largely observe similar empirical results, and
 203 we use these methods mostly as baselines to build intuition.
 204

205 3.1.2 DIRECT PREFERENCE OPTIMIZATION

206
 207 Direct Preference Optimization (DPO) Rafailov et al. (2023) is an offline RL Levine et al. (2020)
 208 alternative to the classical RLHF optimization pipeline. It is a suitable algorithm for agent fine-tuning,
 209 as it can use fully offline data and does not require online rollouts. The original formulation in the
 210 pure text generation setting is based on the RL problem in Eq. 1 and considers feedback of pairwise
 211 comparisons $(\mathbf{h}, \mathbf{a}^w, \mathbf{a}^l)$, where \mathbf{h} is a single prompt and \mathbf{a}^w and \mathbf{a}^l are two responses with $\mathbf{a}^w \succ \mathbf{a}^l$
 212 indicating that \mathbf{a}^w is preferred over \mathbf{a}^l . The DPO objective then minimizes the following loss:
 213

$$214 \mathcal{L}_{\text{DPO}}(\pi_{\theta}; \mathcal{D}) = -\mathbb{E}_{(\mathbf{h}, \mathbf{a}^w, \mathbf{a}^l) \sim \mathcal{D}} \left[\log \sigma \left(\left(\beta \log \frac{\pi_{\theta}(\mathbf{a}^w|\mathbf{h}^w)}{\pi_{\text{ref}}(\mathbf{a}^w|\mathbf{h}^w)} \right) - \left(\beta \log \frac{\pi_{\theta}(\mathbf{a}^l|\mathbf{h}^l)}{\pi_{\text{ref}}(\mathbf{a}^l|\mathbf{h}^l)} \right) \right) \right] \quad (4)$$

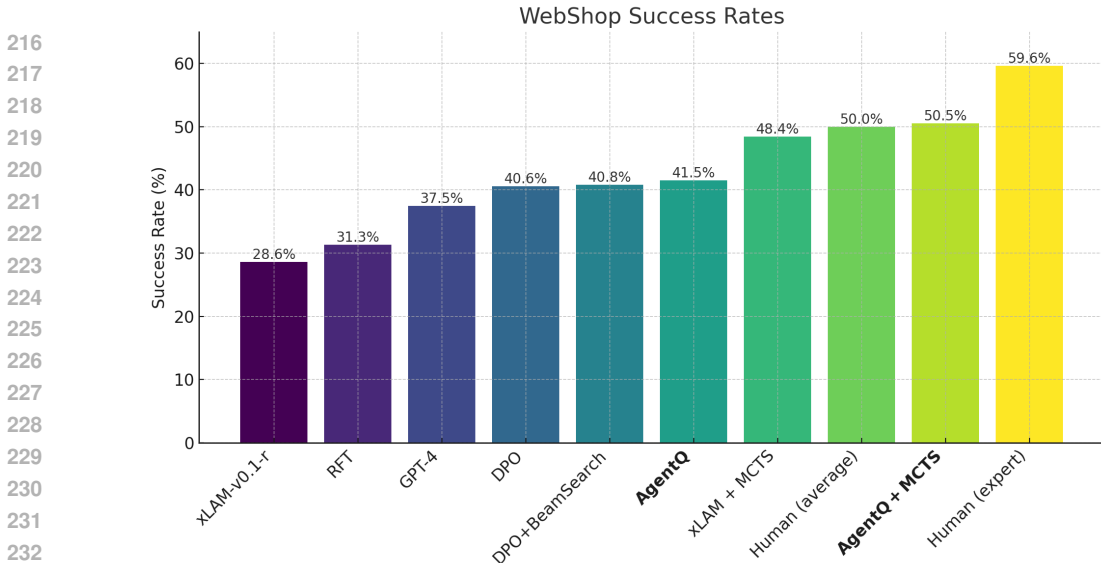


Figure 1: Success rate of different approaches on the WebShop Yao et al. (2022) task. All models are based on xLAM-v0.1-r Zhang et al. (2024c). RFT and DPO over xLAM-v0.1-r demonstrate improvements in performance from 28.6% to 31.3% and 37.5% respectively. However, these methods still lag behind average human performance of 50.0%. Our approach, Agent Q + MCTS achieves a significant gain (76.57% relative improvement) over the base model, outperforming average human performance on WebShop with a success rate of 50.5%.

While the algorithm was developed in a bandit setting Hejna et al. (2024); Rafailov et al. (2024) have extended it to the multi-turn setting in Eq. 2 with preferences over trajectories. In our setting, we can directly utilize this objective as:

$$\mathcal{L}_{\text{T-DPO}}(\pi_{\theta}; \mathcal{D}) = -\mathbb{E}_{(\tau_w, \tau_l) \sim \mathcal{D}} \left[\log \sigma \left(\left(\sum_{t=0}^{|\tau^w|} \beta \log \frac{\pi_{\theta}(\mathbf{a}_t^w | \mathbf{h}_t^w)}{\pi_{\text{ref}}(\mathbf{a}_t^w | \mathbf{h}_t^w)} \right) - \left(\sum_{t=0}^{|\tau^l|} \beta \log \frac{\pi_{\theta}(\mathbf{a}_t^l | \mathbf{h}_t^l)}{\pi_{\text{ref}}(\mathbf{a}_t^l | \mathbf{h}_t^l)} \right) \right) \right] \quad (5)$$

One bottleneck for the practical deployment of the algorithm is the need for a reference model π_{ref} during optimization, which requires more computational resources. Instead in our settings, we slightly modify the algorithm using an off-policy replay buffer, which aggregates trajectory data, as well as likelihoods of the generated actions. During the optimization step, we sample tuples of trajectories and the corresponding likelihoods under the data generation (reference) density, which eliminates the need for a separate reference model.

4 PRELIMINARY APPROACH WITH OUTCOME SUPERVISION

In this section we will outline preliminary experimental results, which will build the base understanding for our further experiments. We use the AgentOhana xLAM-v0.1-r model Zhang et al. (2024c), which is a fine-tune of a pre-trained Mixtral-8x7B-Instruct-v0.1 model Jiang et al. (2024) on a mix of agentic applications, including WebShop SFT data. We also incorporate the same agent configuration¹ specified by the AgentLite Liu et al. (2024) work to ensure a fair comparison between our fine-tuned model and the xLAM base model performance. We evaluate all approaches on the WebShop environment Yao et al. (2022), where the agent needs to find particular products by browsing a simulated web shop. The environment comes with a set of 12,087 pre-defined tasks (corresponding to specific products to find), which we split into a train set of 11,000 tasks, which we use for further agent fine-tuning and a set of 1,087 held-out tasks, which we use for zero-shot evaluation. We show success rates (exact product match) for different approaches in Fig. 1. The base xLAM-v0.1-r model achieves success rate of 28.6% on the test tasks. All other methods are based on

¹<https://github.com/SalesforceAIRResearch/xLAM>

270 outcome-based supervision only, depending on whether a particular attempt was successful or not.
 271 We see that further RFT training, using a STaR-like algorithm [Zelikman et al. \(2022\)](#) on the trajectory
 272 level, as outlined in Sec. 3.1.1, achieves success rate of 31.3%, which is a small improvements of
 273 2.7% over the initial model. This is not surprising since the base model is already trained as an agent
 274 on the environment with supervised fine-tuning on demonstrations. Our next experiment fine-tunes
 275 the base model using the trajectory-level DPO algorithm, as outlined in Eq. 5 in Sec. 3.1.2 using
 276 successful trajectories as preferred over failed ones. This approach also uses only outcome-level
 277 supervision, but unlike the RFT baseline can utilize failed trajectories as well, which improves the
 278 agent performance by 9.3% over RFT agent to 40.6% success rate. We also evaluate this model with
 279 beam search for the action generation, which can be considered a form of planning the horizon of a
 280 single environment action (which still consists of multiple simple actions) [Rafailov et al. \(2023\)](#), but
 281 it only yields marginal improvement over the base model. These findings match results on reasoning
 282 for math problems [Pang et al. \(2024\)](#) and some recent approaches that also apply DPO to agent
 283 applications [Song et al. \(2024\)](#); [Xi et al. \(2024\)](#).

284 Despite the additional reinforcement learning training, our agents are still not able to match the
 285 average human performance on this environment. We identify that one of the core failure modes of
 286 the DPO policy is that it executes a greedy search when looking for matches to the product query.
 287 For example, for every search query, the WebShop environment yields a number of pages of results.
 288 However, we find that the model nearly always greedily searches for the best matching item in the
 289 first page of results rather than using the "[NEXT]" and "[PREV]" buttons to navigate between pages,
 290 essentially deploying a weak exploration strategy.

291 5 AGENT SEARCH

293 As we discovered in the previous section, while training based on outcome supervision with DPO
 294 yields meaningful improvement, the model is still not able to match human performance due to
 295 its limited exploration. In this section we will explore endowing the agent with additional search
 296 capability via MCTS.

298 5.1 MONTE-CARLO TREE SEARCH OVER WEB-PAGES

300 The Monte Carlo Tree Search (MCTS) algorithm [Kocsis and Szepesvári \(2006\)](#) employed in this
 301 work follows closely the one in [Hao et al. \(2023\)](#) and consists of four phases: selection, expansion,
 302 simulation, and backpropagation. Each phase plays a critical role in balancing exploration and
 303 exploitation while iteratively refining the policy.

304 We formulate the web agent execution as tree search over web-pages. The state is represented as
 305 described in Appendix A and consist of the summary of the agent’s history and the DOM tree of
 306 the current web-page. Unlike board games, such as Chess or Go [Silver et al. \(2017b\)](#) the complex
 307 web-agent action space we use is open-format and variable. Instead we will use the base model as an
 308 action-proposal distribution and sample a fixed amount of possible actions at each node (web-page).
 309 Once we select and execute an action in the browser we traverse the next web-page, which together
 310 with the updated history becomes the new node.

311 5.1.1 ACTION SELECTION WITH AI PROCESS SUPERVISION

313 The selection phase uses the Upper Confidence Bound (UCB1) formulation of MCTS also used by
 314 [Hao et al. \(2023\)](#) to select nodes which aims to balance exploration and exploitation. With some
 315 abuse of notation we will also denote the agent state with \mathbf{h}_t . We consider the value function $Q(\mathbf{h}_t, \mathbf{a})$
 316 which represents the estimated value (chance of success) represents the estimated value of taking
 317 action \mathbf{a} in the state \mathbf{h}_t . At each new node \mathbf{h}_t we sample K proposal actions from the base model
 318 $\mathbf{a}_t^1, \dots, \mathbf{a}_t^K$. We initialize all values $Q(\mathbf{h}_t, \mathbf{a}_t^i), i = 1, \dots, K$ to zero. The web-based environment
 319 does not provide intermediate rewards to guide the search, so we incorporate AI-based critique to
 320 provide process supervision at the step level to guide the exploration process. We use the base model
 321 to produce a feedback score for each action by asking it to rank the generated actions by its perceived
 322 utility in helping the agent complete the user task.

323 We query the feedback model for multiple iterations, each time removing the best action selected
 from the previous iteration from the list, until we have a full ranking of all actions. The full AI

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

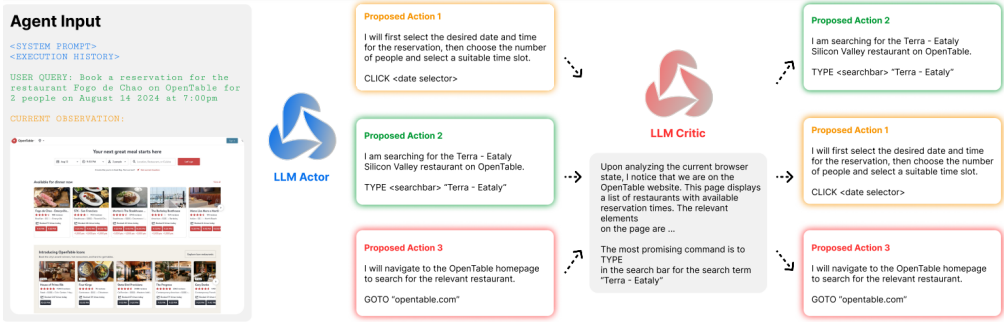


Figure 2: The policy proposes K actions at every step during inference time search. The critic, also initialized as the same base LLM model used by the policy, ranks the actions proposed by the policy. This ranking is used to guide node selection after expansion and used to construct preference pairs during policy training.

feedback process is demonstrated in Figure 2. After the initial selection, we select the actions to explore based on the standard MCTS UCB1 formulation:

$$\mathbf{a}_t^* = \arg \max_{\mathbf{a}_t^1, \dots, \mathbf{a}_t^K} \left[Q(\mathbf{h}_t, \mathbf{a}) + c_{\text{exp}} \cdot \sqrt{\frac{\log N(\mathbf{h}_t)}{1 + N(\mathbf{h}_{t+1})}} \right], \quad (6)$$

where $N(\mathbf{h}_t)$ is the visitation frequency of state \mathbf{h}_t , and c_{exp} is an exploration constant. For each rollout added to the tree, we start at the root node and follow the child states that maximize the UCB1 score until we reach a leaf node. This process is repeated for each tree/prompt in the batch.

5.1.2 EXPANSION AND BACKTRACKING

Based on the preceding section, we select and execute an action in the browser environment to reach a new node (page). Beginning from the selected state node’s trace, we roll out the trajectory using the current policy π_θ until a terminal state is reached. The environment returns a reward at the end of the trajectory, R , where $R = 1$ if the agent was successful and $R = 0$ otherwise. We then backpropagate this reward by updating the values of each node bottom up from the leaf node to the root as follows:

$$\begin{aligned} Q(\mathbf{h}_t, \mathbf{a}_t^i) &\leftarrow \frac{Q(\mathbf{h}_t, \mathbf{a}_t^i)N(\mathbf{h}_t, \mathbf{a}_t^i) + R}{N(\mathbf{h}_t, \mathbf{a}_t^i) + 1} \\ N(\mathbf{h}_t, \mathbf{a}_t^i) &\leftarrow N(\mathbf{h}_t, \mathbf{a}_t^i) + 1 \end{aligned} \quad (7)$$

Each state node tracks two values: $Q(\mathbf{h}_t, \mathbf{a}_t^i)$, the average reward for passing through state \mathbf{h}_t and choosing action \mathbf{a}_t^i , and $N(\mathbf{h}_t, \mathbf{a}_t^i)$, the number of times this state action pair was visited during search (and $N(\mathbf{h}_t) = \sum_{i=1}^K N(\mathbf{h}_t, \mathbf{a}_t^i)$). The backpropagation updates correctly maintain these values.

5.2 IMPROVING ZERO-SHOT PERFORMANCE WITH REINFORCEMENT LEARNING

Training large foundation models with offline Snell et al. (2022) or off-policy Chebotar et al. (2023) reinforcement learning at scale has still remained challenging. At the same time online (on-policy) reinforcement learning Stiennon et al. (2022); Ouyang et al. (2022) is not scalable to real interactive environments. Instead, we follow a line of recent works, which apply the DPO algorithm Rafailov et al. (2023; 2024) at the step level in multi-step reasoning problems in mathematical domains Xie et al. (2024); Hwang et al. (2024); Chen et al. (2024); Lai et al. (2024b); Lu et al. (2024); Setlur et al. (2024b); Zhang et al. (2024f). Our approach is most similar to Xie et al. (2024); Chen et al. (2024); Zhang et al. (2024f) who also use the branching nature of tree search to produce step-level preference pairs. We will also use this approach in our setting due to its simplicity, scalability and prior success in smaller scale (non-interactive) reasoning applications.

We will generate a dataset of preference pairs $\mathcal{P} = \{\mathbf{h}_t, \mathbf{a}_t^w, \mathbf{a}_t^l\}$ where we make sure both actions were explored. We then optimize the DPO objective in Eq. 4 on the node level. We will leverage a

Algorithm 1 MCTS Guided Direct Preference Optimization

378
379
380 **Input:** π_{θ_0} : initial LLM policy, \mathcal{D}_T : dataset of tasks the agent must complete in the environment,
381 N : number of iterations, B : number of samples per iteration, T : MCTS tree depth, \mathcal{B} : replay
382 buffer, $\theta_{\text{threshold}}$: value threshold in (9), K : number of actions to sample for MCTS
383 **Output:** π_{θ_N} , the trained LLM policy
384 **for** $i = 1$ to N **do**
385 $\pi_{\text{ref}} \leftarrow \pi_{\theta_i}, \pi_{\theta_i} \leftarrow \pi_{\theta_{i-1}}$
386 Sample a batch of B tasks from \mathcal{D}_T
387 **for** each task in batch **do**
388 Initialize the root node \mathbf{h}_0
389 **for** $t = 1$ to T **do**
390 **Selection:** Traverse tree from the root node to a leaf node using tree policy (UCB1; 6)
391 **Trajectory Rollout:** From the selected node’s trace, roll out the trajectory using
392 π_{θ_i} until a terminal state is reached
393 **Backpropagation:** Backpropagate the value estimate bottom-up (7)
394 **end for**
395 Collect trajectories from rollouts and store them in replay buffer \mathcal{B}
396 **end for**
397 Construct preference pairs $\mathcal{D}_P = \{(\mathbf{h}_t, \mathbf{a}_t^w, \mathbf{a}_t^l)\}_{t=1}^{T-1}$ where $\mathbf{h}_t \sim \mathcal{D}_P$. For each node at step
398 level t , compare each pair of child nodes, and construct the pair of generated actions ($\mathbf{a}^w, \mathbf{a}^l$) if the
399 values of taking the action, $|Q(\mathbf{h}_t, \mathbf{a}^w) - Q(\mathbf{h}_t, \mathbf{a}^l)| > \theta_{\text{threshold}}$, where $Q(\mathbf{h}_t, \mathbf{a}^w)$ and $Q(\mathbf{h}_t, \mathbf{a}^l)$
400 are computed using (9)
401 Optimize LLM policy π_{θ_i} using DPO objective in Eq. (4) with \mathcal{D}_P and π_{ref}
402 **end for**

403 theoretical result below to guide the construction of these preferences. We can make a number of
404 modifications to Theorem 6.1 from Setlur et al. (2024b) to incorporate the interactive nature of the
405 web environment dynamics to obtain the following result:

406 **Theorem 1.** Consider a policy that optimizes the objective in Eq. 2 on trajectories generated by
407 π_{ref} and that at each node \mathbf{h}_t we have preferences generated accordingly to $p(\mathbf{a}_t^w \succ \mathbf{a}_t^l | \mathbf{h}_t) \propto$
408 $\sigma(Q(\mathbf{h}_t, \mathbf{a}_t^w) - Q(\mathbf{h}_t, \mathbf{a}_t^l))$, then the policy which optimizes the DPO objective in Eq. 4 is identical
409 to the optimal RL policy

$$410 \pi^*(\mathbf{a} | \mathbf{h}_t) \propto \pi_{\text{ref}}(\mathbf{a} | \mathbf{h}_t) \exp(Q(\mathbf{h}_t, \mathbf{a}) / \beta) \quad (8)$$

411 *Proof.* The proof follows directly from the proof of Theorem 6.1 in Setlur et al. (2024b) and the
412 control as inference arguments in Rafailov et al. (2024); Levine (2018). \square

414 That is, we can approximate the optimal RL policy if we generate preferences under the optimal
415 value function (or an approximation thereof). Since the outcome success provides limited supervision
416 we also incorporate process supervision through the AI feedback as outlined in Section 5.1.1. We
417 interpret the ranking of possible actions by the model to be driven by an implicit value function.
418 Similar semantics was used in Koh et al. (2024), where GPT-4 was used as a zero-shot value function,
419 while here we ask the model to instead reason over the given potential actions and provide rankings
420 instead. This self-rewarding approach has shown promise in the RLHF setting Yuan et al. (2024) and
421 we utilize it for our agent setting as well. Under this formulation, we compute the state-action value
422 as an average:

$$423 Q(\mathbf{h}_t, \mathbf{a}_t^i) = \alpha \tilde{Q}(\mathbf{h}_t, \mathbf{a}_t^i) + (1 - \alpha) \hat{Q}(\mathbf{h}_t, \mathbf{a}_t^i) \quad (9)$$

424 where $\tilde{Q}(\mathbf{h}_t, \mathbf{a}_t^i)$ is the empirical value estimated through MCTS backpropagation and $\hat{Q}(\mathbf{h}_t, \mathbf{a}_t^i)$ is
425 a value estimate based on the ranking of the action \mathbf{a}_t^i by the process supervision AI model. We
426 then create preferences over pairs of actions which are above a certain value threshold $|Q(\mathbf{h}_t, \mathbf{a}_t^w) -$
427 $Q(\mathbf{h}_t, \mathbf{a}_t^l)| \geq \theta_{\text{threshold}}$. The full outline of our RL approach is shown in Algorithm 1.

429 5.3 FULL WEBSHOP RESULTS

430 The full range of results and baselines is shown in Figure 1. We see that equipping the agent with
431 search capabilities at test time significantly boost success rates from 28.6% to 48.4% when using

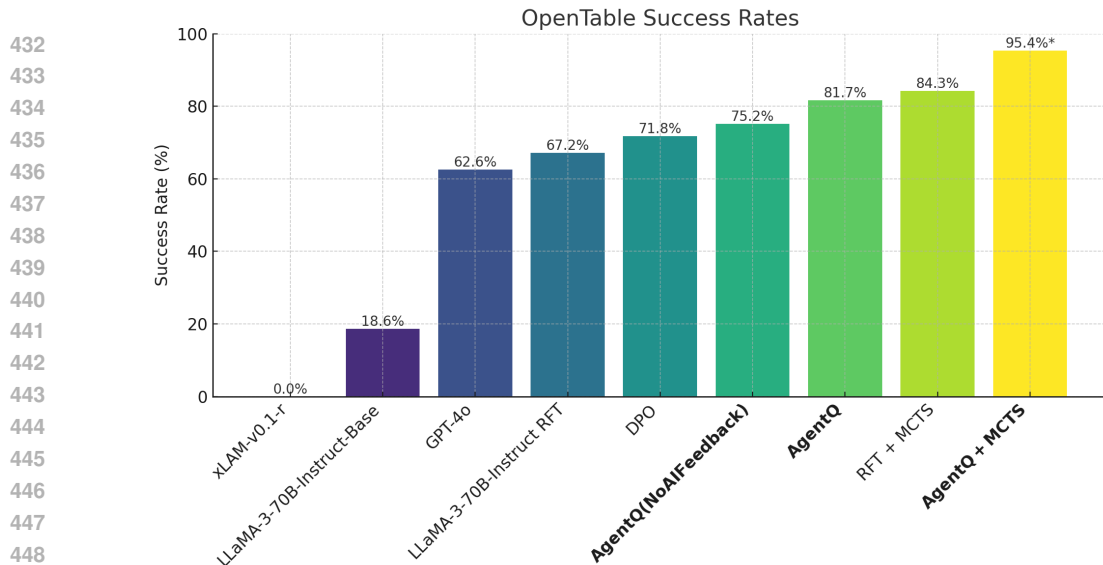


Figure 3: Success rate of different approaches on OpenTable. All models unless otherwise stated are based on LLaMA-3-70B-Instruct [Touvron et al. \(2023\)](#). Using DPO and RFT with MCTS further improves performance from 18.6% to 71.8% and 84.3% respectively. We show that Agent Q in itself achieves 81.7% and Agent Q + MCTS significantly outperforms all other techniques, with a performance of **95.4%** on OpenTable.

MCTS on top of the base xLAM-v0.1-r model, approaching close to the average human performance of 50.0% and significantly out-performing the zero-shot performance of the DPO model trained with outcome supervision. We further fine-tune the base model using the approach outlined in Algorithm 1, which yields an improvement of 0.9% over the base DPO model. Using MCTS on top of the trained Agent Q model further improves performance to 50.5% slightly out-performing the average human success rates. We find that the ability to search at test time is a significant paradigm shift from zero-shot agents, even with significant RL training. Furthermore, while dense-level supervision improves over purely outcome-based one, the improvement is modest on WebShop. This is because the environment requires relatively short trajectories, and the model is capable to learn credit assignment purely from outcome supervision. We will further explore more complex real world environment, which requires longer-range credit assignment.

6 SCALING TO REAL WORLD WEBSITES

In this section we will investigate scaling the Agent Q framework to real use cases on live websites, in particular bookings on OpenTable. We carried out initial experiments with the xLAM-v0.1-r model, which proved to weak for the task achieving an initial success rate of 0.0%. Instead we shifted to the LLaMa 70B Instruct model, which was able to achieve some non-trivial initial success. For description of our real-world environment, consult Appendix B.

The base xLAM-v0.1-r model achieves a success rate of 0.0%, largely from failing to follow instructions for the general web navigation instructions used for live websites, contrary to the simplified observation and action space used in WebShop. We instead initialize the base policy with the LLaMa-3 70B Instruct model, which achieves a zero-shot success rate of 18.6%. We do a single round of RFT on 600 successful trajectories which improves the success rate to 67.2% already out-performing the the GPT-4o model zero-shot performance with a success rate of 62.6%. For all other baselines we adopt the RFT model as the reference policy, due to the relatively low success rate of original LLaMa 3 70B Instruct model.

In this environment, training with outcome-supervision only DPO further improves performance by 4.6% to 71.8% but significantly under-performs the full Agent Q pipeline which achieves a zero-shot success rate of 81.7% We hypothesizes that this is due to the fact that OpenTable is a significantly more challenging environment, which requires almost twice as many steps to complete as WebShop,

486 so the agent benefits from fine-grained supervision and credit assignment. We further ablate the role
487 of the intermediate AI feedback process supervision during training as outlined in Eq. 9 and use
488 MCTS with online Q values computed from outcome rewards only. This setting still outperforms
489 training with trajectory-level DPO (75.2% versus 71.8%) likely due to the more fine-grained credit
490 assignment that the branching tree search provides to the agent. However, zero-shot performance
491 is still meaningfully worse than using intermediate process-level supervision and the full Agent Q
492 achieves 6.5% higher success rate at 81.7%.

493 Similar to the WebShop experiment we see a step level increase in capability from allowing the
494 model to search at inference time, with the base RFT model achieving 84.3% success with MCTS,
495 outperforming the Agent Q zero-shot performance of 81.7% success. However, if we carry out
496 additional MCTS search using the Agent Q model as the base policy we achieve a significant 95.4%
497 success rate.

499 7 DISCUSSION

501 In this work we developed algorithms for autonomous improvement of web-agents with limited
502 human supervision. While most prior works build frameworks around existing models without
503 additional training, we specifically seek to fine-tune pre-trained models for web navigation tasks
504 based on synthetic reasoning and search data. While we achieve significant improvement in model
505 capabilities on our target domain, many research questions remain.

506 **Design of reasoning algorithms.** The core challenge for our web agents is the weak reasoning
507 capabilities, which limit the agent’s exploration and search strategy. In our approach we used process-
508 level supervision from a separate critic model, which we prompt to rank possible agent actions. This
509 is in contrast to works in mathematical reasoning where PRMs are usually trained to classify the
510 correctness of individual steps [Lightman et al. \(2023\)](#), while other agent works [Koh et al. \(2024\)](#)
511 have prompted models as zero-shot value functions. Furthermore, while we spent significant effort in
512 training the agent policy, we maintain a frozen critic, which would likely also benefit from additional
513 fine-tuning. We defer exploration of these design choices to further work.

514 **Choice of search algorithm.** We used MCTS search due to the approach’s prior success in mathemat-
515 ical and code reasoning tasks. However, agent models executing MCTS on live environments might
516 require significant number of risky interactions and a different search strategy might be more suitable.
517 Recent works such as [Lehnert et al. \(2024\)](#); [Gandhi et al. \(2024\)](#) have even suggested directly learning
518 to optimally search and explore in reasoning tasks using meta-reinforcement learning. We believe
519 this is a promising research direction for autonomous agents, which we will pursue in further work.

520 **Discrepancy between zero-shot vs search results.** Similar to some recent works that focus on code
521 and reasoning, we observe significant gap between zero-shot agent performance and performance of
522 the agent equipped with search capabilities [Snell et al. \(2024\)](#); [Brown et al. \(2024\)](#). Investigating
523 these trade-offs at scale and the potential effect of different search/optimization approaches.

524 **Online safety and interaction.** The design of agent Q allows for largely autonomous exploration,
525 self-evaluation and improvement with limited human intervention. However, the agent might make a
526 significant number of mistakes in its search process which might be difficult to fix/reverse, especially
527 for safety-critical online transactions, such as communications/email, payments, filings etc. This
528 limits the scope of websites that Agent Q can be safely deployed and we might require additional
529 safety critics and human-in-the-loop training setups.

REFERENCES

- 540
541
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
543 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
544 *arXiv preprint arXiv:2303.08774*, 2023.
- 545 Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan
546 Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: A family of highly capable
547 multimodal models. *arXiv preprint arXiv:2312.11805*, 1, 2023.
- 548 Anthropic. Introducing the next generation of claude, 2024. URL
549 [IntroducingthenextgenerationofClaude](#).
- 550
551 Hao Bai, Yifei Zhou, Mert Cemri, Jiayi Pan, Alane Suhr, Sergey Levine, and Aviral Kumar. Digirl:
552 Training in-the-wild device-control agents with autonomous reinforcement learning, 2024.
553
- 554 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones,
555 Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson,
556 Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson,
557 Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile
558 Lukosuute, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado,
559 Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec,
560 Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom
561 Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei,
562 Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional ai: Harmlessness
563 from ai feedback, 2022.
- 564 Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi,
565 Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. Graph
566 of thoughts: Solving elaborate problems with large language models. *Proceedings of the AAAI*
567 *Conference on Artificial Intelligence*, 38(16):17682–17690, March 2024. ISSN 2159-5399. doi:
568 10.1609/aaai.v38i16.29720. URL <http://dx.doi.org/10.1609/aaai.v38i16.29720>.
- 569 Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and
570 Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling,
571 2024. URL <https://arxiv.org/abs/2407.21787>.
- 572
573 Noam Brown and Tuomas Sandholm. Superhuman ai for multiplayer poker. *Science*, 2019.
- 574 Yevgen Chebotar, Quan Vuong, Alex Irpan, Karol Hausman, Fei Xia, Yao Lu, Aviral Kumar, Tianhe
575 Yu, Alexander Herzog, Karl Pertsch, Keerthana Gopalakrishnan, Julian Ibarz, Ofir Nachum,
576 Sumedh Sontakke, Grecia Salazar, Huong T Tran, Jodilyn Peralta, Clayton Tan, Deeksha Manju-
577 nath, Jaspiar Singht, Brianna Zitkovich, Tomas Jackson, Kanishka Rao, Chelsea Finn, and Sergey
578 Levine. Q-transformer: Scalable offline reinforcement learning via autoregressive q-functions,
579 2023.
- 580 Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai Fan. Step-level value preference optimization for
581 mathematical reasoning, 2024. URL <https://arxiv.org/abs/2406.10858>.
- 582
583 Wei Chen and Zhiyuan Li. Octopus v2: On-device language model for super agent, 2024.
584
- 585 Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and
586 Yu Su. Mind2web: Towards a generalist agent for the web. In *NeurIPS Datasets and Benchmarks*
587 *Track*, 2023. URL <https://openreview.net/forum?id=kiYqb03wqw>.
- 588 Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin,
589 Percy Liang, and Tatsunori B. Hashimoto. AlpacaFarm: A simulation framework for methods that
590 learn from human feedback, 2024.
- 591
592 Kanishk Gandhi, Denise Lee, Gabriel Grand, Muxin Liu, Winson Cheng, Archit Sharma, and
593 Noah D. Goodman. Stream of search (sos): Learning to search in language, 2024. URL <https://arxiv.org/abs/2404.03683>.

- 594 Jonathan Gray, Adam Lerer, Anton Bakhtin, and Noam Brown. Human-level performance in no-press
595 diplomacy via equilibrium search, 2021.
596
- 597 Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek
598 Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, Wolfgang Macherey, Arnaud
599 Doucet, Orhan Firat, and Nando de Freitas. Reinforced self-training (rest) for language modeling,
600 2023.
- 601 Izzeddin Gur, Hiroki Furuta, Austin V Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and
602 Aleksandra Faust. A real-world webagent with planning, long context understanding, and program
603 synthesis. In *ICLR*, 2024. URL <https://openreview.net/forum?id=9JQtrumvg8>.
604
- 605 Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu.
606 Reasoning with language model is planning with world model, 2023.
- 607 Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan,
608 and Dong Yu. Webvoyager: Building an end-to-end web agent with large multimodal models.
609 *ArXiv*, 2024. URL <https://api.semanticscholar.org/CorpusID:267211622>.
- 610 Joey Hejna, Rafael Rafailov, Harshit Sikchi, Chelsea Finn, Scott Niekum, W. Bradley Knox, and Dorsa
611 Sadigh. Contrastive preference learning: Learning from human feedback without reinforcement
612 learning. In *The Twelfth International Conference on Learning Representations*, 2024. URL
613 <https://openreview.net/forum?id=iX1RjVQ0Dj>.
614
- 615 Samuel Holt, Max Ruiz Luyten, and Mihaela van der Schaar. L2mac: Large language model
616 automatic computer for extensive code generation, 2024.
- 617 Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazhen Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan
618 Wang, Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent: A visual language model for gui agents.
619 *ArXiv*, 2023.
620
- 621 Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot
622 planners: Extracting actionable knowledge for embodied agents, 2022a.
- 623 Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan
624 Tompson, Igor Mordatch, Yevgen Chebotar, Pierre Sermanet, Noah Brown, Tomas Jackson, Linda
625 Luu, Sergey Levine, Karol Hausman, and Brian Ichter. Inner monologue: Embodied reasoning
626 through planning with language models, 2022b. URL <https://arxiv.org/abs/2207.05608>.
627
- 628 Jan Humplik, Alexandre Galashov, Leonard Hasenclever, Pedro A. Ortega, Yee Whye Teh, and
629 Nicolas Heess. Meta reinforcement learning as task inference, 2019. URL <https://arxiv.org/abs/1905.06424>.
630
- 631 Hyeonbin Hwang, Doyoung Kim, Seungone Kim, Seonghyeon Ye, and Minjoon Seo. Self-explore to
632 avoid the pit: Improving the reasoning capabilities of language models with fine-grained rewards,
633 2024.
- 634 Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris
635 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand,
636 Gianna Lengyel, Guillaume Bour, Guillaume Lample, L elio Renard Lavaud, Lucile Saulnier, Marie-
637 Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le
638 Scao, Th eophile Gervet, Thibaut Lavril, Thomas Wang, Timoth ee Lacroix, and William El Sayed.
639 Mixtral of experts, 2024.
640
- 641 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
642 Narasimhan. Swe-bench: Can language models resolve real-world github issues?, 2024.
- 643 Barret Zoph John Schulman, Christina Kim, Jacob Hilton, Jacob Menick, Jiayi Weng, Juan Fe-
644 lipe Ceron Uribe, Liam Fedus, Michael Pokornyy Luke Metz, Rapha Gontijo Lopes, Shengjia Zhao,
645 Arun Vijayvergiya, Eric Sigler, Adam Perelman, Chelsea Voss, Mike Heaton, Joel Parish, Dave
646 Cummings, Rajeev Nayak, Valerie Balcom, David Schnurr, Tomer Kaftan, Chris Hallacy, Nicholas
647 Turley, Noah Deutsch, Vik Goel, Jonathan Ward, Aris Konstantinidis, Wojciech Zaremba, Long
Ouyang, Leonard Bogdonoff, Joshua Gross, David Medina, Sarah Yoo, Teddy Lee, Ryan Lowe,

- 648 Dan Mossing, Joost Huizinga, Roger Jiang, Carroll Wainwright and Diogo Almeida, Steph Lin,
649 Marvin Zhang, Kai Xiao, Katarina Slama, Steven Bills, Alex Gray, Jan Leike, Jakub Pachocki, Phil
650 Tillet, Shantanu Jain, Greg Brockman, Nick Ryder, Alex Paino, Qiming Yuan, Clemens Winter,
651 Ben Wang, Mo Bavarian, Igor Babuschkin, Szymon Sidor, Ingmar Kanitscheider, Mikhail Pavlov,
652 Matthias Plappert, Nik Tezak, Heewoo Jun, William Zhuk, Vitchyr Pong, Lukasz Kaiser, Jerry
653 Tworek, Andrew Carr, Lilian Weng, Sandhini Agarwal, Karl Cobbe, Vineet Kosaraju, Alethea
654 Power, Stanislas Polu, Jesse Han, Raul Puri, Shawn Jain, Benjamin Chess, Christian Gibson,
655 Oleg Boiko, Emy Parparita, Amin Tootoonchian, Kyle Kosic, and Christopher Hesse. Introducing
656 chatgpt, 2022. URL <https://openai.com/blog/chatgpt#OpenAI>.
- 657 Jikun Kang, Xin Zhe Li, Xi Chen, Amirreza Kazemi, Qianyi Sun, Boxing Chen, Dong Li, Xu He,
658 Quan He, Feng Wen, Jianye Hao, and Jun Yao. Mindstar: Enhancing math reasoning in pre-trained
659 llms at inference time, 2024. URL <https://arxiv.org/abs/2405.16265>.
- 660 Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *Machine Learning:
661 ECML 2006: 17th European Conference on Machine Learning Berlin, Germany, September 18-22,
662 2006 Proceedings*, pages 282–293. Springer, 2006.
- 663
664 Jing Yu Koh, Stephen McAleer, Daniel Fried, and Ruslan Salakhutdinov. Tree search for language
665 agent models, 2024. URL <https://jykoh.com/search-agents/paper.pdf>.
- 666 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
667 language models are zero-shot reasoners, 2022.
- 668 Hanyu Lai, Xiao Liu, Iat Long Iong, Shuntian Yao, Yuxuan Chen, Pengbo Shen, Hao Yu, Hanchen
669 Zhang, Xiaohan Zhang, Yuxiao Dong, and Jie Tang. Autowebglm: Bootstrap and reinforce a large
670 language model-based web navigating agent, 2024a.
- 671
672 Xin Lai, Zhuotao Tian, Yukang Chen, Senqiao Yang, Xiangru Peng, and Jiaya Jia. Step-dpo: Step-
673 wise preference optimization for long-chain reasoning of llms, 2024b. URL [https://arxiv.org/
674 abs/2406.18629](https://arxiv.org/abs/2406.18629).
- 675 Lucas Lehnert, Sainbayar Sukhbaatar, DiJia Su, Qinqing Zheng, Paul Mcvay, Michael Rabbat, and
676 Yuandong Tian. Beyond a*: Better planning with transformers via search dynamics bootstrapping,
677 2024. URL <https://arxiv.org/abs/2402.14083>.
- 678
679 Sergey Levine. Reinforcement learning and control as probabilistic inference: Tutorial and review,
680 2018.
- 681 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,
682 review, and perspectives on open problems, 2020.
- 683
684 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike,
685 John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step, 2023.
- 686
687 Zhiwei Liu, Weiran Yao, Jianguo Zhang, Le Xue, Shelby Heinecke, Rithesh Murthy, Yihao Feng,
688 Zeyuan Chen, Juan Carlos Niebles, Devansh Arpit, Ran Xu, Phil Mui, Huan Wang, Caiming Xiong,
689 and Silvio Savarese. Bolaa: Benchmarking and orchestrating llm-augmented autonomous agents,
690 2023.
- 691
692 Zhiwei Liu, Weiran Yao, Jianguo Zhang, Liangwei Yang, Zuxin Liu, Juntao Tan, Prafulla K. Choubey,
693 Tian Lan, Jason Wu, Huan Wang, Shelby Heinecke, Caiming Xiong, and Silvio Savarese. Agentlite:
A lightweight library for building and advancing task-oriented llm agent system, 2024.
- 694
695 Zimu Lu, Aojun Zhou, Ke Wang, Houxing Ren, Weikang Shi, Juntong Pan, Mingjie Zhan, and
696 Hongsheng Li. Step-controlled dpo: Leveraging stepwise error for enhanced mathematical
reasoning, 2024. URL <https://arxiv.org/abs/2407.00782>.
- 697
698 Tula Masterman, Sandi Besen, Mason Sawtell, and Alex Chao. The landscape of emerging ai agent
699 architectures for reasoning, planning, and tool calling: A survey, 2024.
- 700
701 Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher
Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted
question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.

- 702 Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher
703 Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou,
704 Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. Webgpt:
705 Browser-assisted question-answering with human feedback, 2022.
706
- 707 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
708 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser
709 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan
710 Leike, and Ryan Lowe. Training language models to follow instructions with human feed-
711 back. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Ad-
712 vances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Asso-
713 ciates, Inc., 2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/file/
714 b1efde53be364a73914f58805a001731-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/b1efde53be364a73914f58805a001731-Paper-Conference.pdf).
- 715 Jiayi Pan, Yichi Zhang, Nicholas Tomlin, Yifei Zhou, Sergey Levine, and Alane Suhr. Autonomous
716 evaluation and refinement of digital agents, 2024.
- 717 Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho, He He, Sainbayar Sukhbaatar, and Jason
718 Weston. Iterative reasoning preference optimization, 2024.
719
- 720 Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen, Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei
721 Huang, and Huajun Chen. Reasoning with language model prompting: A survey, 2023.
722
- 723 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
724 Finn. Direct preference optimization: Your language model is secretly a reward model. In *Thirty-
725 seventh Conference on Neural Information Processing Systems*, 2023. URL [https://arxiv.org/
726 abs/2305.18290](https://arxiv.org/abs/2305.18290).
- 727 Rafael Rafailov, Joey Hejna, Ryan Park, and Chelsea Finn. From r to q^* : Your language model is
728 secretly a q-function, 2024.
- 729 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
730 optimization algorithms, 2017.
731
- 732 Amrith Setlur, Saurabh Garg, Xinyang Geng, Naman Garg, Virginia Smith, and Aviral Kumar. RL on
733 incorrect synthetic data scales the efficiency of llm math reasoning by eight-fold, 2024a.
734
- 735 Amrith Setlur, Saurabh Garg, Xinyang Geng, Naman Garg, Virginia Smith, and Aviral Kumar. RL
736 on incorrect synthetic data scales the efficiency of llm math reasoning by eight-fold, 2024b. URL
737 <https://arxiv.org/abs/2406.14532>.
- 738 David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez,
739 Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy Lillicrap, Fan
740 Hui, Laurent Sifre, George van den Driessche, Thore Graepel, and Demis Hassabis. Mastering the
741 game of go with deep neural networks and tree search. *Nature*, 2017a.
- 742 David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez,
743 Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without
744 human knowledge. *Nature*, 550(7676):354–359, 2017b.
745
- 746 Avi Singh, John D. Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, Xavier Garcia, Peter J.
747 Liu, James Harrison, Jaehoon Lee, Kelvin Xu, Aaron Parisi, Abhishek Kumar, Alex Alemi, Alex
748 Rizkowsky, Azade Nova, Ben Adlam, Bernd Bohnet, Gamaleldin Elsayed, Hanie Sedghi, Igor
749 Mordatch, Isabelle Simpson, Izzeddin Gur, Jasper Snoek, Jeffrey Pennington, Jiri Hron, Kathleen
750 Kenealy, Kevin Swersky, Kshiteej Mahajan, Laura Culp, Lechao Xiao, Maxwell L. Bileschi,
751 Noah Constant, Roman Novak, Rosanne Liu, Tris Warkentin, Yundi Qian, Yamini Bansal, Ethan
752 Dyer, Behnam Neyshabur, Jascha Sohl-Dickstein, and Noah Fiedel. Beyond human data: Scaling
753 self-training for problem-solving with language models, 2024.
- 754 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally
755 can be more effective than scaling model parameters, 2024. URL [https://arxiv.org/abs/2408.
03314](https://arxiv.org/abs/2408.03314).

- 756 Charlie Victor Snell, Ilya Kostrikov, Yi Su, Sherry Yang, and Sergey Levine. Offline rl for natural
757 language generation with implicit language q learning. In *The Eleventh International Conference*
758 *on Learning Representations*, 2022.
- 759 Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. Trial and error:
760 Exploration-based trajectory optimization for llm agents, 2024.
- 761 Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,
762 Dario Amodei, and Paul Christiano. Learning to summarize from human feedback, 2022.
- 763 Theodore R. Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L. Griffiths. Cognitive architec-
764 tures for language agents, 2024.
- 765 Fahim Tajwar, Anikait Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie, Stefano
766 Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of llms should leverage suboptimal,
767 on-policy data, 2024.
- 768 Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Haitao Mi, and Dong Yu. Toward
769 self-improvement of llms via imagination, searching, and criticizing, 2024.
- 770 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
771 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, Dan Bikel, Lukas Blecher, Cris-
772 tian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu,
773 Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,
774 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel
775 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,
776 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,
777 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,
778 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh
779 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen
780 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,
781 Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models,
782 2023.
- 783 Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia
784 Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process- and
785 outcome-based feedback, 2022.
- 786 Junyang Wang, Haiyang Xu, Jiabo Ye, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao
787 Sang. Mobile-agent: Autonomous multi-modal mobile device agent with visual perception. *arXiv*,
788 2024a.
- 789 Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang
790 Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations, 2024b.
- 791 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
792 ury, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models,
793 2023.
- 794 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le,
795 and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. *Neural*
796 *Information Processing Systems*, 2022.
- 797 Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao
798 Yu, and Lingpeng Kong. Os-copilot: Towards generalist computer agents with self-improvement.
799 *arXiv*, 2024. URL <https://api.semanticscholar.org/CorpusID:265149992>.
- 800 Zhiheng Xi, Yiwen Ding, Wenxiang Chen, Boyang Hong, Honglin Guo, Junzhe Wang, Dingwen
801 Yang, Chenyang Liao, Xin Guo, Wei He, Songyang Gao, Lu Chen, Rui Zheng, Yicheng Zou,
802 Tao Gui, Qi Zhang, Xipeng Qiu, Xuanjing Huang, Zuxuan Wu, and Yu-Gang Jiang. Agentgym:
803 Evolving large language model-based agents across diverse environments, 2024. URL <https://arxiv.org/abs/2406.04151>.

- 810 Yuxi Xie, Anirudh Goyal, Wenyue Zheng, Min-Yen Kan, Timothy P. Lillicrap, Kenji Kawaguchi, and
811 Michael Shieh. Monte carlo tree search boosts reasoning via iterative preference learning, 2024.
812
- 813 John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan,
814 and Ofir Press. Swe-agent: Agent-computer interfaces enable automated software engineering,
815 2024.
- 816 Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable
817 real-world web interaction with grounded language agents, 2022.
818
- 819 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik R.
820 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In *NeurIPS*,
821 2023a. URL <https://openreview.net/forum?id=5Xc1ecx01h>.
- 822 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
823 React: Synergizing reasoning and acting in language models, 2023b.
824
- 825 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and
826 Jason Weston. Self-rewarding language models, 2024.
827
- 828 Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Keming Lu, Chuanqi Tan, Chang Zhou,
829 and Jingren Zhou. Scaling relationship on learning mathematical reasoning with large language
830 models, 2023.
- 831 Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. Star: Bootstrapping reasoning with
832 reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
833
- 834 Yuexiang Zhai, Hao Bai, Zipeng Lin, Jiayi Pan, Shengbang Tong, Yifei Zhou, Alane Suhr, Saining
835 Xie, Yann LeCun, Yi Ma, and Sergey Levine. Fine-tuning large vision-language models as
836 decision-making agents via reinforcement learning, 2024.
- 837 Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang, Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qingwei
838 Lin, and Saravan Rajmohan. Ufo: A ui-focused agent for windows os interaction. *arXiv*, 2024a.
839 URL <https://api.semanticscholar.org/CorpusID:267211622>.
- 840 Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu.
841 Appagent: Multimodal agents as smartphone users, 2023.
842
- 843 Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang
844 Yu. Appagent: Multimodal agents as smartphone users. *arXiv*, 2024b. URL <https://api.semanticscholar.org/CorpusID:262053313>.
- 845 Jianguo Zhang, Tian Lan, Rithesh Murthy, Zhiwei Liu, Weiran Yao, Juntao Tan, Thai Hoang,
846 Liangwei Yang, Yihao Feng, Zuxin Liu, Tulika Awalgaonkar, Juan Carlos Niebles, Silvio Savarese,
847 Shelby Heinecke, Huan Wang, and Caiming Xiong. Agentohana: Design unified data and training
848 pipeline for effective agent learning, 2024c.
- 849 Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi Jin. Codeagent: Enhancing code generation with
850 tool-integrated agent systems for real-world repo-level coding challenges, 2024d.
- 851 Xuan Zhang, Chao Du, Tianyu Pang, Qian Liu, Wei Gao, and Min Lin. Chain of preference
852 optimization: Improving chain-of-thought reasoning in llms, 2024e.
853
- 854 Xuan Zhang, Chao Du, Tianyu Pang, Qian Liu, Wei Gao, and Min Lin. Chain of preference
855 optimization: Improving chain-of-thought reasoning in llms, 2024f. URL <https://arxiv.org/abs/2406.09136>.
- 856 Zhuosheng Zhang and Aston Zhang. You only look at screens: Multimodal chain-of-action agents.
857 *ArXiv*, 2023. URL <https://api.semanticscholar.org/CorpusID:262053313>.
- 858 Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Language
859 agent tree search unifies reasoning acting and planning in language models, 2024a.
860
861
862
863

864 Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,
865 Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic
866 web environment for building autonomous agents. In *ICLR*, 2024b.
867
868 Yifei Zhou, Andrea Zanette, Jiayi Pan, Sergey Levine, and Aviral Kumar. Archer: Training language
869 model agents via hierarchical multi-turn rl, 2024c.
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

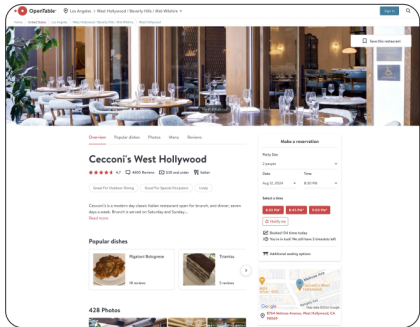
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Agent Input

```
<SYSTEM PROMPT>
<EXECUTION HISTORY>

USER QUERY: Book a reservation for the
restaurant Ceconni's on OpenTable for
2 people on June 17 2024 at 7:00pm

CURRENT OBSERVATION:
```



Agent Output

```
PLAN:
1. Select the date (May 22, 2024) for the
   reservation.
2. Choose the time (7:00 PM) for the
   reservation.
3. Select the party size (4 people) for the
   reservation.
4. Click on the "Find a table" button to book
   the reservation.

THOUGHT: I am currently on the Ceconni's
restaurant page on OpenTable, and I need to
select the date and time for the reservation
and choose the number of guests. I will focus
on selecting the date, time, and party size
for the reservation.

COMMANDS: CLICK <select>Date</select>

STATUS: CONTINUE
```

Figure 4: We provide the following input format to the Agent, consisting of the system prompt, execution history, the current observation as a DOM representation, and the user query containing the goal. We divide our Agent output format into an overall step-by-step plan, thought, a command, and a status code.

A AGENT FORMULATION

We consider a general POMDP setup $(\mathcal{O}, \mathcal{S}, \mathcal{A}, T, R, \mu_0, \gamma)$ where \mathcal{O} denotes the observation space, \mathcal{S} the unobserved state space, \mathcal{A} the action space, $T(s_{t+1}|s_t, \mathbf{a}_t)$ the transition distribution (in this case the dynamics of a web browser), $R(s, \mathbf{a})$ the reward function (in this work we use sparse rewards of 1/0 representing success/failure), $\mu_0(s_0)$ the initial state distribution, and γ the discount factor, which we set to 1. A POMDP is the most suitable framework to model web interactions for several reasons - first novel environments, which the agent is unfamiliar with require exploration in order to locate the task objective, consistent with the meta-reinforcement learning as task inference view [Humplik et al. \(2019\)](#). Moreover, the real web is dynamic, which creates partial observability of the current state each time the agent is deployed - i.e. it does not a priori know current booking availability before attempting to do it. We will outline the main parts of our web agent below.

The agent observation $\mathbf{o}_t \in \mathcal{O}$ are commands/information given by the user and the web browser. The first observation \mathbf{o}_1 is a user text instruction, such as

”Book reservation for restaurant Ceconni’s on OpenTable for 4 people on May 22 2024 at 7:00 PM”

for example and a browser home page. Subsequent observations consist of web pages from the browser, represented as a HTML DOM format. Occasionally for some tasks the agent might ask for confirmation/feedback from the user, which then also becomes part of the observation.

The agent actions $\mathbf{a}_t \in \mathcal{A}$ are composite, based on agent history \mathbf{h}_t . Our base approach is a ReAct agent [Yao et al. \(2023b\)](#) with a preliminary planning step (PlanReAct) [Liu et al. \(2023\)](#) with few additional components.

- **Planning** For the first action after the initial observation we leverage the base LLM’s planning capabilities [Huang et al. \(2022a\)](#) and prompt the agent to generate a plan $\mathbf{a}_1^{\text{plan}} \sim \pi(\mathbf{a}_1^{\text{plan}}|\mathbf{h}_1)$ of sequential steps to execute in language.
- **Reasoning** Subsequently all actions consist of a thought action $\mathbf{a}_t^{\text{tht}} \sim \pi(\mathbf{a}_t^{\text{tht}}|\mathbf{h}_t)$, which is reasoning step [Wei et al. \(2022\)](#).

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

- **Environment action** Next we generate the browser interaction command $\mathbf{a}_t^{\text{env}} \sim \pi(\mathbf{a}_t^{\text{env}}|\mathbf{h}_t, \mathbf{a}_t^{\text{tht}})$, which consists of a finite set of options like "CLICK [ELEMENT ID]", "SCROLL", "TYPE [CONTENT]" or "ASK USER [CONTENT]" etc.. This is the only part of the action generation, which interacts with the environment.
- **Explanation action** After the environment interaction action has been generated, we additionally prompt the model for an explanation action $\mathbf{a}_t^{\text{expl}} \sim \pi(\mathbf{a}_t^{\text{expl}}|\mathbf{h}_t, \mathbf{a}_t^{\text{tht}}, \mathbf{a}_t^{\text{env}})$.

We denote the step action \mathbf{a}_t as a tuple of plan, thought, environment and explanation actions for the first step and thought, environment and explanation actions for subsequent steps. When optimizing models we consider the joint likelihood

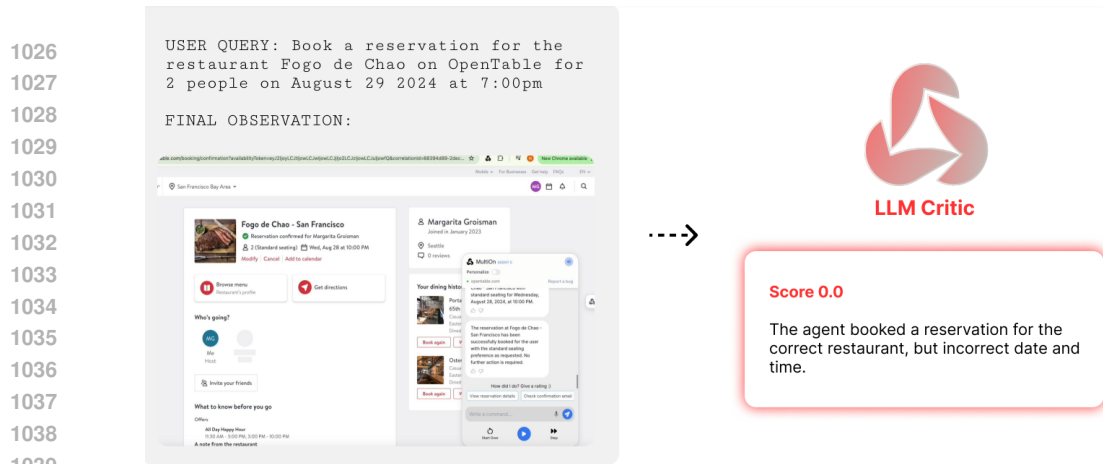
$$\log \pi(\mathbf{a}_1|\mathbf{h}_1) = \log \pi(\mathbf{a}_1^{\text{expl}}|\mathbf{h}_1, \mathbf{a}_1^{\text{env}}, \mathbf{a}_1^{\text{tht}}, \mathbf{a}_1^{\text{plan}}) + \log \pi(\mathbf{a}_1^{\text{env}}|\mathbf{h}_1, \mathbf{a}_1^{\text{tht}}, \mathbf{a}_1^{\text{plan}}) + \log \pi(\mathbf{a}_1^{\text{tht}}|\mathbf{h}_1, \mathbf{a}_1^{\text{plan}}) + \log \pi(\mathbf{a}_1^{\text{plan}}|\mathbf{h}_1) \tag{10}$$

for the initial action and

$$\log \pi(\mathbf{a}_t|\mathbf{h}_t) = \log \pi(\mathbf{a}_t^{\text{expl}}|\mathbf{h}_t, \mathbf{a}_t^{\text{env}}, \mathbf{a}_t^{\text{tht}}) + \log \pi(\mathbf{a}_t^{\text{env}}|\mathbf{h}_t, \mathbf{a}_t^{\text{tht}}) + \log \pi(\mathbf{a}_t^{\text{tht}}|\mathbf{h}_t)$$

for subsequent actions, unlike some prior works [Zhai et al. \(2024\)](#), which down-weight the reasoning likelihood.

The agent state is the current state of the web, which may not be observable. In this POMDP formulation we also need to build an agent memory component \mathbf{h}_t . Prior works have used the entire trajectory of observations and actions, however HTML DOMs can be hundred of thousands of tokens long. Moreover realistic web-tasks can require many more interactions than static benchmarks such as WebShop [Yao et al. \(2022\)](#) and WebArena [Zhou et al. \(2024b\)](#), which most prior works use. This makes it impractical to use full web trajectories due to limited context windows, potential out-of-distribution issues and practical inference speed and cost. Instead, we build the history representation of the agent as $\mathbf{h}_t = (\mathbf{a}_1, \dots, \mathbf{a}_{t-1}, \mathbf{o}_t)$. That is, the agent history consists of the actions generated so far and the current browser state. With some abuse of notation we will also refer to this as the agent state. Even though only the environment action is used for interacting with the browser, we construct the agent thought and explanation actions to act as a form of inner monologue [Huang et al. \(2022b\)](#) and adequately represent its state and intentions. This allows us to use a significantly more compact history representation. We should note that, while only the environment action affects the browser state, the planning, reasoning and explanation components affect subsequent decisions due to conditioning. Because of this reason, when we optimize the agent, we compute likelihoods over the composite action.



1041 Figure 5: At the end of a trajectory, a GPT-4-V evaluator is called to provide feedback on the agent’s
 1042 performance given the final observation and action history to determine the success score. The model
 1043 is prompted with a condensed execution history of the trajectory and the screenshot of the final state.
 1044 The success metric is a binary 0/1 value.

1046 B OPENTABLE ENVIRONMENT

1048 In OpenTable, the agent is tasked with booking a restaurant reservation for a user. The agent must
 1049 find a restaurant page on the OpenTable site, look for a reservation at a certain date and time, choose
 1050 seating options that align with a user’s preference and submit the user contact information to complete
 1051 the task successfully. Since OpenTable is a live environment and is difficult to programmatically
 1052 measure metrics for, we use a language model, GPT-4-V to collect rewards for each trajectory,
 1053 based on the following metrics: (1) date and time set correctly, (2) party size set correctly, (3) user
 1054 information entered correctly, and (4) clicked complete reservation. The task is marked as completed
 1055 if each of the above constraints are satisfied. The outcome supervision setup is shown in Figure
 1056 5. We experimented with using LLaMa 70B for outcome supervision as well, but discovered that
 1057 vision capabilities significantly improve the success classification accuracy (as measured by human
 1058 validation). At the time of writing no open source vision-language model of sufficient capability was
 1059 available, hence we opted to use GPT-4-V. We believe that as more open-source multi-modal models
 1060 become available we can switch to a fully self-supervised pipeline.

1061 To generate queries for the OpenTable benchmark dataset, we programmatically generate a diverse set
 1062 of user queries by combining the restaurant name, desired date and time, and user information.

1063 Navigating on live websites pose a wide variety of challenges. For example, consider that the user
 1064 specifies a restaurant in a different city than the location the browser is initialized in, the model will
 1065 have to take extra steps to find the restaurant. Further, if the exact user requested date and time are
 1066 not available, the model may have to choose the closest available reservation slot. Lastly, if there are
 1067 preferences, such as indoor or outdoor seating options that the model is presented with, the desired
 1068 behavior is to interact with the user to determine the best course of action. OpenTable presents a
 1069 complex set of challenges for web navigation agents; the number of steps required to complete the
 1070 task is on average 13.9 steps, over double the average number of steps for Webshop, 6.8.

1071 For the observation space for this environment, we design an intermediate state representation that
 1072 crawls the raw HTML content of a website to retrieve relevant visual components, and highlight
 1073 interactive elements to the model. The agent is allowed the actions, "CLICK [ID]", "GOTO [URL]",
 1074 "TYPE [ID] [TEXT]", "SUBMIT [ID]", "CLEAR [ID]", "SCROLL [UP/DOWN]", and "ASK USER
 1075 HELP". For OpenTable experiments, we use the LLaMA-3-70B-Instruct model as the initial policy.
 1076 We find that the superior reasoning abilities of this class of model is required for effective task
 1077 completion, which is necessary to produce the diverse success and failure trajectories required to
 1078 effectively improve the policy.

1079