PROPOSER-AGENT-EVALUATOR (PAE): AUTONOMOUS SKILL DISCOVERY FOR FOUNDATION MODEL INTERNET AGENTS

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

034

037

Paper under double-blind review

ABSTRACT

The vision of a broadly capable and goal-directed agent, such as an Internetbrowsing agent in the digital world and a household humanoid in the physical world, has rapidly advanced, thanks to the generalization capability of foundation models. Such a generalist agent needs to have a large and diverse skill repertoire, such as finding directions between two travel locations and buying specific items from the Internet. If each skill needs to be specified manually through a fixed set of human-annotated instructions, the agent's skill repertoire will necessarily be limited due to the quantity and diversity of human-annotated instructions. In this work, we address this challenge by proposing Proposer-Agent-Evaluator(PAE), a complete working system that enables foundation model agents to autonomously discover and practice skills in the wild. At the heart of PAE is a context-aware task proposer that autonomously proposes tasks for the agent to practice with context information of the websites such as user demos or even just the name of the website itself. Then, the agent policy attempts those tasks with thoughts and actual web operations in the real world with resulting trajectories evaluated by an autonomous model-based success evaluator. The success evaluation serves as the reward signal for the agent to refine its policies through RL. We validate PAE on challenging vision-based web navigation, using both real-world and self-hosted websites from WebVoyager (He et al., 2024) and WebArena (Zhou et al., 2024a). Our results show that PAE significantly improves the zero-shot generalization capability of VLM Internet agents (more than 30% relative improvement) to both unseen tasks and websites. Our model also achieves an absolute advantage of over 10% (from 22.6% to 33.0%) comparing to other state-of-the-art open source VLM agents including Qwen2VL-72B. To the best of our knowledge, this work represents the first working system to apply autonomous task proposal with RL for agents that generalizes real-world human-annotated benchmarks with sota performances. We plan to release our models and code to facilitate further research.

039 1 INTRODUCTION

040 The vision of broadly capable and goal-directed agent, such as an Internet-browsing agent in the 041 digital world and a household humanoid in the physical world, has long captured our imagination. 042 With recent advancements in foundation models (OpenAI, 2024; GeminiTeam, 2024), this vision is 043 no longer a distant dream. These developments have significantly accelerated the progress of gener-044 alist agents (Liu et al., 2023b) in real-world decision-making scenarios such as navigating through online websites to make travel plans (He et al., 2024) and solving real-user Github issues (Jimenez et al., 2024), making them a rapidly emerging research frontier. To succeed in these decision-making 046 domains, goal-directed post-training is often needed to elicit long-horizon reward-maximizing be-047 haviors such as information seeking (Hong et al., 2023a) and recovery from mistakes (Bai et al., 048 2024), instead of only imitating the most probable actions in the pre-training corpus. 049

A crucial requirement for a successful post-training approach is to endow the generalist agent with
 a large and diverse goal-directed skill repertoire. This can include finding directions between two
 travel locations and buying specific items from the Internet, which the agent can then exploit to solve
 real-world tasks proposed by users. However, manually specifying the skills (Deng et al., 2023) (i.e.
 through a static set of human-annotated instruction templates such as "Find the driving directions

054 and estimated time to travel from Location A to Location B") will likely result in a limited skill 055 repertoire. First of all, generating high-quality human-annotated task templates can be expensive, making it impractical to scale up. The use of a small set of task templates fails to capture the range 057 of skills an agent needs for the full breadth of the real world, leading to distribution shift problems 058 when deployed at the test time. Furthermore, human-generated instructions have limited diversity due to human creativity (Wang et al., 2023b), failing to capture the long-tail distribution of realworld tasks that the agent needs to solve. With these disadvantages, it naturally raises the research 060 question: instead of requiring users to manually define tasks for foundation model agents, can these 061 agents automatically discover and practice potentially useful skills on their own? 062

063 In order to discover its own skills and im-064 prove autonomously, such an agent would need to be able to propose semantically 065 meaningful tasks and then determine if 066 it was successful in performing them. 067 Such success detections can then serve 068 as reward signals to apply Reinforcement 069 Learning (RL) to optimize the agents. While prior works have explored the use 071 of foundation models to propose skills 072 to the agents to practice and detecting 073 successes in simplified environments such 074 as games (Du et al., 2023; Colas et al., 075 2023) and robotics with limited number of



Figure 1: An overview of our method showing the main components of our autonomous skill discovery framework, that endows the agent with autonomously discovered skills to prepare for future human requests.

scenes (Zhang et al., 2023b; Zhou et al., 2024b), little has been understood in terms of whether
such diverse skills can generalize to real-world human request such as web agents and what the key
design decisions are to improve such generalization.

079 To this end, our main contribution is to propose a fully working system, Proposer-Agent-Evaluator(PAE), for foundation model agents (in particular, Internet agents) to autonomously dis-081 cover new skills without any human supervision and such new skills can be effectively exploited to solve unseen real-world human-annotated tasks in a zero-shot manner. In this way, the training workflow can be easily scaled to make use of diverse self-generated instructions in large quantities to 083 enrich the skill repertoire of the agent. PAE is built with the awareness the asymmetric capabilities 084 of sota VLMs as task proposers/ evaluators and as agents (discussed in Section 6) in some realistic 085 task settings such as web agents and designed to make best use of this asymmetry. Intuitively, VLMs 086 are very good at confirming whether a specific product has been added to the shopping cart (e.g. by 087 looking at the final screenshot to see if the shopping cart contains the product), while less good at 880 actually navigating the web to find the product and add it to the cart. To obtain the most robust re-089 ward signal without accessing the hidden state information, we apply an image-based evaluator that 090 only provides sparse 0/1 rewards based on the final outcome. To propose feasible and realistic tasks, 091 PAE employs a class of context-aware task proposers where the context of functions and constraints 092 crucially define what actions are supported by the specific environments (e.g., creating a reddit post) while others may not be supported (e.g., checking the protected information of other users). Such 093 context can be implicitly defined from different sources and are shown to be effective, such as user 094 demos and even website name alone! Finally, we design an additional reasoning step before the 095 agent outputs actual actions, which enables the agents to better reflect on its skills and results in a 096 significant improvement in its generalization capability to unseen human-annotated tasks. 097

098 The scope of our experiments covers challenging end-to-end vision-based web navigation, where the observation space simply contains the screenshot of the current web page and the action space contains primitive web operations such as clicking on links and typing into text boxes. We validate the 100 effectiveness of PAE framework with realistic web-navigation benchmarks, including 16 domains 101 both from online websites like Amazon from WebVoyager (He et al., 2024) and self-hosted web-102 sites like PostMill from WebArena (Zhou et al., 2024a). In our experiments, we find that PAE with 103 LLaVa-1.6 (Liu et al., 2024) as the agent policy can autonomously discover useful skills through 104 interactions with various websites without any human supervisions. More importantly, our results 105 demonstrate that these skills can zero-shot transfer to unseen test instructions and even unseen test 106 websites. On websites from WebVoyager and WebArena, PAE attains a 30% relative improvement 107 in average success rate, enabling LLaVa-1.6-7B to achieve performance comparable with LLaVa1.6-34B fine-tuned with demonstration data despite using 5x fewer test-time compute. Compared to other state-of-the-art open-sourced VLM agents, including Qwen2VL-72B (Yang et al., 2024a), our model achieves an absolute performance gain of over 10% (22.6% to 33.0%). To the best of our knowledge, this work is the first to develop a working system of autonomous skill discovery for foundation model agents that directly generalizes to real-world human-annotated benchmarks.

114 2 RELATED WORKS

113

115 Foundation model agents. Thanks to the generalization capabilities of Large Language Models 116 (LLMs) (Brown et al., 2020; Llama3Team, 2024; GeminiTeam, 2024) and Vision Language Models (VLMs) (OpenAI, 2024; Liu et al., 2024; Wang et al., 2024b; Liu et al., 2023a), recent works have 117 successfully extended such agents to more general real-world use cases (Bai et al., 2024; Zheng et al., 118 2024; He et al., 2024; Zhang et al., 2023a; Zhou et al., 2024a; Koh et al., 2024; Gur et al., 2021; 119 Furuta et al., 2024). Besides constructing prompting wrappers around proprietary VLMs (Zhang 120 et al., 2023a; He et al., 2024; Zheng et al., 2024; Xie et al., 2024; Yang et al., 2024b; Wang et al., 121 2023a) and fine-tuning open-source VLMs with expert demonstrations (Gur et al., 2021; Hong et al., 122 2023b; Furuta et al., 2024; Zhang & Zhang, 2024; Zeng et al., 2023; Chen et al., 2023), a recent trend 123 has emerged involving the interactive improvement of LLM/VLM, in particular web/GUI agents, 124 through autonomous evaluator feedback (Pan et al., 2024; Bai et al., 2024; Putta et al., 2024), where 125 evaluator LLMs/VLMs are prompted to evaluate the success of the agents to serve as the reward 126 signal. This approach aims to elicit goal-oriented and reward-optimizing behaviors from foundation 127 models with minimal human supervision. However, these methods still depend on a static set of human-curated task templates, which can constrain their potential and scalability. Our work intro-128 duces a novel framework where agents can *discover* and practice the skills they find useful, thereby 129 eliminating the reliance on predefined and human-curated task templates. This approach opens up 130 new possibilities for scalability and adaptability in training autonomous LLM/VLM agents. 131

132 Self-generated instructions. Self-generated instructions for improving LLMs have been shown to 133 be effective in single-turn LLM alignment (Wang et al., 2023b; Yuan et al., 2024; Wu et al., 2024; Wang et al., 2024a) and reasoning (Pang et al., 2024) domains without interactions with an external 134 environment. AgentGen (Hu et al., 2024) employs a similar methodology to fine-tune LLM agents 135 with ground-truth trajectories in self-generated environments and tasks. However, its feasibility 136 in the self-play agent setting with RL and autonomous evaluators has not been understood. The 137 closest works to ours employ autonomous RL and foundation model task proposers to simplified 138 environments such as games (Zhang et al., 2024a; Faldor et al., 2024; Colas et al., 2023; 2020) and 139 robotics settings with limited number of scenes (Du et al., 2023; Zhang et al., 2023b; Zhou et al., 140 2024b). While they have shown that the use of autonomous RL combined with foundation model 141 task proposers can help the agent learn diverse skills, this work takes an important step forward to 142 study when those skills can generalize to human requests in realistic benchmarks in the context of 143 web agents and what the best design choices are for such generalization.

144 Unsupervised skill discovery in deep RL. Unsupervised skill discovery has been an important 145 research direction in the field of traditional deep RL literatrue (Achiam et al., 2018; Eysenbach 146 et al., 2018) where various algorithms have been developed to discover new robotic skills such as 147 humanoid walking without the need of explicitly defined reward functions. Common algorithms in 148 this field aim to discover every possible skill (either meaningful skills like walking or less meaningful 149 ones like random twisting) through either maximizing the mutual information between different 150 states and skill latent vectors (Campos et al., 2020; Laskin et al., 2022; Sharma et al., 2020), or maximizing the divergence of each skill as measured in a metric space (Park et al., 2022; 2023; 151 2024). In contrast, our work only discovers meaningful skills as specified through natural language 152 instructions with the help of pre-trained foundation models, significantly reducing the search space 153 of skills in LLM/VLM agent applications with complex state spaces. 154

- 155
- 156 157

3 PROPOSER-AGENT-EVALUATOR (PAE): AUTONOMOUS SKILL DISCOVERY SYSTEM FOR FOUNDATION MODEL AGENTS

Next, we will explain the technical contributions of this paper. In this section, we will define the
general system of PAE including a task proposer, an agent policy, and an autonomous evaluator.
We will begin by formalizing the learning goal of this system and detailing the roles of each key
component in the system. Then we will walk through our practical algorithm in the system. In the
section to follow, we will provide the example of applying PAE to VLM Internet agents.

162 Problem setup We begin by formalizing the problem setup of autonomous skill discovery for real-163 world agents. The learning goal of PAE is to find a reward-maximizing policy π parameterized by θ 164 in a contextual Markov Decision Process (MDP) environment defined by $\mathcal{M} = \{S, \mathcal{A}, \mathcal{T}, \mathcal{R}, H, C\}$, 165 where \mathcal{S}, \mathcal{A} are the state space and action space respectively, and H is the horizon within which 166 the agent must complete the task. We assume that the agent has access to the environment and can collect online roll-out trajectories through accessing the dynamics model \mathcal{T} as a function of 167 determining the next states given the current states and actions. Crucially, we assume that the 168 ground-truth task distribution C and the reward function $\mathcal R$ are hidden during training and we have to use a proxy task distribution \hat{C} and reward function $\hat{\mathcal{R}}$ instead. Consider the setting of training a 170 real-world Internet agent. The dynamics model $\mathcal T$ would be a simulated browser environment that 171 the Internet agent can interact with. The ground-truth task distribution C might be the distribution 172 of tasks that would be asked by the real users when the Internet agent is deployed and a possible 173 choice for the reward function \mathcal{R} might be whether the agent has satisfactorily completed the tasks 174 for the real users. In such a real-world setting, although the agent can freely access resources from 175 the Internet through a simulated browser environment during training, assuming knowledge of the 176 ground-truth task distribution and reward function is impractical. Therefore, we employ VLM-based 177 task proposers \hat{C} and reward model $\hat{\mathcal{R}}$ as proxies. The desired outcome is that improving the policy 178 π_{θ} with $\hat{\mathcal{C}}$ and $\hat{\mathcal{R}}$ can lead to an improved policy that can successfully generalize to the ground-truth 179 task distribution and reward functions which are only used as evaluations.

Key components Figure 1 shows the interplay between the key components in our framework, 181 including a context-aware task proposer, an agent policy, and an autonomous evaluator. The role 182 of the **task proposer** $\hat{\mathcal{C}}$ is to serve as a proxy to improve on the ground-truth task distribution \mathcal{C} 183 during the learning process. However, it might be unrealistic to expect the task proposer to generate feasible tasks without knowledge of the environment. To provide more context of the functions and 185 constraints of the environment, we assume access to some key information of the environment $z_{\mathcal{M}}$ based on which the tasks $\mathcal{C}(z_{\mathcal{M}})$ are proposed. In the Internet agent example, this key information 187 can be screenshots of the websites from user demos, or even just the name of the website itself if it 188 is a well-known website such as Amazon.com. Similarly, the **autonomous evaluator** \mathcal{R} serves as a 189 proxy of the ground-truth reward function \mathcal{R} . The input to the autonomous evaluator is the current 190 state, the current action from the agent policy, and current task that the agent is attempting. In 191 principle, any RL algorithm can be used to update the **agent policy** π using a dataset \mathcal{D} that stores 192 all the autonomous interaction data. In practice, we instantiate VLM-based task proposers and autonomous evaluators by prompting foundation models and they are kept unchanged throughout 193 our practical algorithm. 194

195 196

4 PROPOSER-AGENT-EVALUATOR FOR VLM INTERNET AGENTS

With the general framework set up, we are now
ready to discuss the concrete instantiation of
PAE in the setting of VLM Internet agents. We
start by introducing the environment of visionbased web navigation and then explain how we
implement the key components from the PAE
in this setting.

- 4.1 VISION-BASED
- 205 WEB BROWSING ENVIRONMENT

We consider the general vision-based web browsing environment (He et al., 2024; Koh et al., 2024). The goal for VLM agents in this environment is to navigate through realistic web pages to complete some user tasks c_t such as "Investigate in the Hugging Face documentation how to utilize the 'Trainer' API for training a model on a custom dataset, and note



Figure 2: An illustration of the observation space and action space of our vision-based web navigation environment. The observation space is augmented with set-of-marks that label each interactable element with a unique number. At each step, the web agent first chooses an element to interact with by referring to its number and then choose the action type to perform on this element (e.g., click, type, and etc.).

the configurable parameters of the Trainer class". As illustrated in Figure 2, each **observation** s_t from the observation space contains only the screenshot of the last web page just like how humans interact with the Internet. To provide better action grounding, we follow the practice from prior works (Zheng et al., 2024; He et al., 2024) to augment the observation space with number marks on 216 top of each interactive element such as web links and text boxes. To execute a web browsing action, 217 the Internet agent can directly output the number of the element to interact with and the correspond-218 ing action such as clicking and typing, without the need of locating the coordinates of each web 219 element. Therefore each web **action** a_t contains the type of the action to perform and the number of 220 the element to interact with. Each episode finishes either when the agent chooses to finish through the "Answer" action or when a maximum number of 10 steps have been reached. In our experiments, we use ground-truth success detectors (based on either human annotations or functional verifiers) 222 and human annotated tasks from WebArena (Zhou et al., 2024a) and WebVoyager (He et al., 2024) 223 to evaluate the performance of different policies. Crucially, both the ground-truth success detector 224 and the distribution of human tasks are kept hidden, which challenges the generalization capability 225 of the learnt skills to generalize to a hidden reward function and task distributions. 226

227 4.2 CONTEXT-AWARE TASK PROPOSER

228 In order to generate a diverse set of feasible tasks, we frame task proposing \hat{C} as a conditional auto-regressive generation based on the context information of the websites. Thanks to the vast 229 pre-training knowledge of relevant context for popular websites like Amazon.com, we find it suffice 230 to use only website name as $z_{\mathcal{M}}$. However, for less common or access restricted websites such as 231 self-hosted websites in WebArena, it is necessary to supply the task proposer with richer context. 232 In the cases of **user demos** being available, we consider an alternative to sample some additional 233 screenshots from the user demos to serve as the context information. In our experiments, we consider 234 both using proprietary models such as Claude-3-Sonnet (Anthropic, 2024) and open-source models 235 such as Qwen2VL-7B Yang et al. (2024a) for the task proposers, with promptsd in Appendix C.

236

237 4.3 IMAGE-BASED OUTCOME EVALUATOR

To take full advantage of the asymmetric capability of SOTA VLMs as agents and as evaluators 238 (experiment results presented in Section 6, we empirically find it reliable for the autonomous evalu-239 ators to complete the easiest evaluation: evaluating the success of the final outcome (Bai et al., 2024; 240 He et al., 2024) based on the final three screenshots and the agents' final answers to provide only 241 0/1 response in the end. Other alternatives such as code-based (Zhang et al., 2024a) or step-based 242 evaluations (Pan et al., 2024) are either impractical without access to hidden state information or too 243 noisy because of the hallucination issues present even in SOTA VLMs. In our experiments, we also 244 consider both using proprietary models such as Claude-3-Sonnet (Anthropic, 2024) and open-source 245 models such as Qwen2VL-7B Yang et al. (2024a), with prompts presented in Appendix C.

4.4 CHAIN-OF-THOUGHT AGENT POLICY

247 Crucially, as the ultimate goal for the agent policy is to complete human requests, the agent should 248 not only learn diverse skills on the proposed tasks but also reflect on the skills learnt so that they 249 can be helpful for unseen human requests. Therefore, we incorporate an additional reasoning step 250 to outputs the agent's chain-of-thought before the actual web operation. This reasoning step is optimized with the RL algorithm just like the actual web operation. Because of the 0/1 reward 251 structure and infrastructure complexity of thousands of distributed fully-functioning web browsers, 252 we employ the most simple online policy optimization algorithm Filtered Behavior Cloning (Filtered 253 BC), that simply imitates all thoughts and actions in successful trajectories with the negative log-254 liklihood loss. We find that this simple policy optimization objective can already lead to an superior 255 generalization capability of the learnt agent. In our experiments, our agent policy is initialized from 256 LLaVa-1.6-Mistral-7B and LLaVa-1.6-Yi-34B (Liu et al., 2024). 257

258 5 EXPERIMENTS

The goal of our experiments is to understand the effectiveness of PAE to complete real-world visual web tasks. Specifically, we design experiments to answer the following questions: (1) Can our autonomous skill discovery framework successfully discover skills useful for zero-shot transfer to tasks from an evaluation task distribution unseen to the task proposer? (2) How does the models trained with PAE compare with other open-source VLM agents? (3) How does the effectiveness of PAE scale with the size and performance of the base model? (4) How does the use of different contexts (e.g. website names and user demos) affect the performance?

266 5.1 Environments

WebVoyager (He et al., 2024) contains a set of 643 tasks spanning 15 websites in the real world such
 as ESPN and Arxiv. As tasks in Google Flights and Booking domain are no longer feasible due to
 website updates, we use the subset of 557 tasks spanning the other 13 websites. Human annotations are carried out for evaluating the success of each trajectory as the ground-truth performance measure.

270			Allrecipes	Amazon	Apple	ArXiv	GitHub	ESPN	Coursera
271	Proprietary	Claude 3.5 Sonnet Claude 3 Sonnet	50.0 15.9	68.3 46.3	60.4 51.1	46.5 39.5	58.5 41.4	27.3 11.3	78.6 45.2
272		Qwen2VL-7B Owen2VL-72B	0.0	0.0	2.3	2.3	0.0	0.0	2.3 48 5
273	Open-source	InternVL2.5-8B	0	0	0	0	0	0	0
274		LLaVa-34B	0	0	2.3	0	2.4	0	0
275	Ours	LLaVa-7B SF1 LLaVa-7B PAE LLaVa-7B PAE (Owen7B)	4.5	37.5	17.5	19.0	4.9 14.6 15.6	0.0	33.3
276	OWS	LLaVa-34B SFT LLaVa-34B PAE	6.8 22.7	26.8 53.7	23.3 38.5	16.3 25.6	4.9 14.6	8.6 13.6	26.8 42.9
277			Cambridge Dictionary	BBC News	Google Map	Google Search	HuggingFace	Wolfram Alpha	Average
278	Proprietary	Claude 3.5 Sonnet Claude 3 Sonnet	86.0 79.1	36.6 40.5	58.5 41.5	30.2 41.9	44.2 37.2	66.7 61.9	50.5 42.4
279		Qwen2VL-7B Qwen2VL-72B	2.3 60.6	0.0 12.5	0.0 16.1	4.7 21.2	0.0 9.1	4.8 36.4	1.4 22.6
280	Open-source	InternVL2.5-8B LLaVa-7B LLaVa-34P	0	0	0	0 0 2.3	2.3 0 2.3	0 0	0.2
281		LLava-54B LLaVa-7B SFT LLaVa-7B PAE	41.9 52.4	7.1	19.5 22.5	9.3 23.3	0	11.9 24.4	14.9
282	Ours	LLaVa-7B PAE (Qwen7B)	62.5 67.4	12.5 16.7	12.9	3.0		36.4 38.1	21.7
283		LLaVa-34B PAE	74.4	39.0	22.0	18.6	25.6	42.9	33.0

284 Table 1: Success rate comparisons on WebVoyager. The results are automatically annotated by Claude 285 Sonnet 3 and human alignment is reported in Figure 4. For PAE, a running average of the evaluation results at 286 each iteration is reported. The final column is a weighted average by the number of tasks on different websites. The results may be different from reported in other papers due to the dynamic nature of online websites. 287

288 WebArena (Zhou et al., 2024a) is a sand-boxed environment that kept an archived version of 5 pop-289 ular websites from different domains, including OpenStreetMap, GitLab, PostMill, a store content management system (CMS), and an E-commerce website (OneStopMarket). It includes in total 812 290 hand-written tasks with functional verifications as the ground-truth reward function. Since GitLab 291 and CMU do not support multi-thread data collection necessary for RL fine-tuning, our experiments 292 are carried out using the task subsets on OpenStreetMap, PostMill, and OneStopMarket. As open-293 source VLM agents fail to achieve non-trivial performances on PostMill and OneStopMarket (Zhou et al., 2024a), we hand rewrote tasks in those two websites and supplement them with verification 295 functions. Due to these practical constraints, the resulting WebArena Easy contains 108 original 296 tasks on OpenStreetMap and 50 rewritten tasks on PostMill and OneStopMarket each. 297

298 5.2 **BASELINE COMPARISONS**

299 We validate the effectiveness of PAE by comparing it with (1) 300 proprietary VLMs, (2) state-of-the-301 art open-source VLMs, and (3) 302 an alterative supervised fine-tuning 303 (SFT) approach. We consider 304 Claude 3 Sonnet and Claude 3.5 305 Sonnet (Anthropic, 2024) for pro-306 prietary VLMs, and Qwen2VL-307 7B, Qwen2VL-72B (Yang et al., 308 2024a), InternVL-2.5-XComposer-309 **7B** (Zhang et al., 2024b), and LLaVa-Next-7B/34B (Liu et al., 310 2024) for SOTA open-source VLMs. 311 All models are prompted similar to 312

		OpenStreetMap	PostMill	OneStopMarket	Average
D	Claude 3.5 Sonnet	38.3	70.0	53.0	50.1
Frophetary	Claude 3 Sonnet	24.3	55.8	41.7	36.0
	Qwen2VL-7B	0.7	10.2	20.2	7.5
	Qwen2VL-72B	16.0	32.8	32.7	23.9
0	InternVL2.5-8B	1.8	0.5	6.0	2.5
Open-source	LLaVa-7B	0.0	0.0	0.0	0.0
	LLaVa-34B	0.9	0.0	0.0	0.5
	LLaVa-7B SFT	15.2	16.8	25.4	18.0
	LLaVa-7B PAE	19.5	21.1	42.3	24.6
Ours	LLaVa-7B PAE (Qwen72B)	17.9	30.6	39.2	26.0
	LLaVa-7B PAE (Qwen7B)	20.2	25.0	28.6	23.1
	LLaVa-7B PAE (User Demos)	21.7	21.5	42.1	25.7

Table 2: Success rate comparisons on WebArena Easy. Success and failure are detected with ground-truth verification functions. For PAE, a running average of the evaluation results at each iteration is reported. The final "Average" column is a weighted average by the number of tasks on different websites.

He et al. (2024) using set-of-marks augmented screenshot observations and including chain-of-313 thought in the action outputs. The prompts are included in Appendix C. As SOTA open-source 314 models struggle to achieve non-trivial performance in the challenging web navigation benchmarks 315 except the largest Qwen2VL-72B, we include another baseline LLaVa-SFT that fine-tunes LLaVa 316 with Claude 3 Sonnet (Anthropic, 2024) agent trajectories on self-generated tasks on 85 real-world 317 websites not included in WebVoyager and WebArena. More details in the data generation for SFT 318 can be found in Appendix E. To study the effects of different contexts for our task proposer, we 319 compare the performance of two variants from PAE as discussed in Section 4.2. LLaVa-34B PAE 320 and LLaVa-7B PAE uses only the name of the website as the context, while LLaVa-7B-PAE (User 321 **Demos**) uses 10 additional screenshots per website from human collected user demos. To investigate the effect of SOTA VLMs on the improvements of PAE, we include two additional baselines 322 LLaVa-34B PAE (Qwen7B) and LLaVa-34B PAE 72B where we use Qwen2VL-7B/Qwen2VL-323 72B as both task proposers and evaluators.

324 5.3 MAIN RESULTS

We present our main baseline comparisons of PAE with other baselines in Table 1, 2, and 3. Overall, 326 comparing to the SFT checkpoint using demonstration data, LLaVa-7B PAE can achieve an average 327 of 7.4% and 10.8% absolute improvement in terms of success rates on WebVoyager and WebArena 328 Easy respectively. A similar improvement of 10.4% on WebVoyager is observed for LLaVa-34B 329 PAE as well, indicating a favorable scaling performance of PAE. As a result, our resulting model 330 LLaVa-34B PAE achieves an absolute success rate of 10.4% on WebVoyaer over the prior state-of-331 the-art open-source VLM agents. Similarly, LLaVa-7B PAE also establishes a new state-of-the-art 332 performance on WebArena Easy, surpassing the prior best performing model Qwen2VL-72B with 333 $10 \times$ more parameters. More importantly, our analysis shows that PAE can enable Internet agents to 334 learn general web browsing capabilities that zero-shot transfer to unseen websites.

How does existing open-source and proprietary mod-

els perform in vision-based web navigation? First, we 337 note the difficulty and significance of real-world vision-338 based web navigation, even for state-of-the-art medium-339 size open-source VLM agents such as Qwen2VL-7B and 340 InternVL2.5-8B with set-of-marks augmented observa-341 tions and chain-of-thought prompting. In particular, on the WebVoyager benchmark, among open-source VLM 342 agents, only the largest Qwen2VL-72B can achieve a 343 non-trivial average success rate of 22.6% on WebVoy-344 ager, while all other open-source agents completely fail 345 on this benchmark with average success rate under 2%. 346 On the other hand, closed-source proprietary models start 347 to show promise in becoming a generalist Internet agent 348 with Claude 3.5 Sonnet achieving an average success rate 349 at 50.5% and 50.1% on WebVoyager and WebArena Easy.

Model	Seen Websites	Unseen Websites
Claude 3 Sonnet	42.4	25.0
Qwen2VL-7B	1.4	1.4
Qwen2VL-72B	22.6	8.3
LLaVa-7B SFT	14.9	9.1
LLaVa-7B PAE	22.3	16.3
LLaVa-7B PAE (Qwen7B)	21.7	13.7
LLaVa-34B SFT	22.2	16.1
LLaVa-34B PAE	33.0	21.4

Table 3: Task success rate comparisons on unseen websites that PAE never interacts with. We select 85 unseen real-world online websites and generate 500 synthetic tasks similar to the procedure in WebVoyager (He et al., 2024). Seen websites are 13 online websites in WebVoyager. Results show that PAE can discover general web browsing skills useful for unseen websites.

Comparing LLaVa-7B SFT and LLaVa-7B, we find that supervised fine-tuning on demonstration
data can significantly improve the general web browsing capabilities of open-source VLM agents.
Even if the SFT demonstration data is collected on out-of-distribution online websites, the general
web browsing capabilities can zero-shot transfer to WebVoyager websites, resulting in a performance
improvement from 0% to 14.9%.

355 Is PAE able to autonomously discover and practice skills useful for unseen evaluation instruc-356 tions? On top of the performance gain from downstream fine-tuning, LLaVa-7B PAE additionally 357 improves the success rate by more than 30% relatively (14.9% to 22.3% on WebVoyager and 18.0% 358 to 24.6% on WebArena Easy). In particular, LLaVa-7B PAE beats LLaVa-7B SFT across the board with substantial improvements on 10 out of 13 websites from WebVoyager and all 3 websites from 359 WebArena Easy, showing the robustness of the PAE framework. In fact, LLaVa-7B PAE even beats 360 the LLaVa-34B SFT (22.3% compared to 22.2%), a model more than 5x larger (7B and 34B), result-361 ing a better performing model with 5x less test-time compute. The release of our models marks a 362 significant advancement of screenshot-based web browsing capabilities of open-source VLM agents 363 from the prior SOTA of 22.6% to 33.0% on WebVoyager. It also enables medium-size VLMs such 364 as LLaVa-7B to beat the prior SOTA Qwen2VL-72B with $10 \times$ more parameters on WebArena Easy. Notably, all of the improvements from PAE are achieved in a self-play setting without any human 366 intervention, only knowing the names of the websites! 367

Is the improvement of PAE bounded by the performance of the evaluator and task proposer 368 **model?** To understand whether the improvement of PAE is bounded by the performance of the 369 model used as the task proposer and evaluator, we replicate the experiments of PAE using open-370 source VLMs (Qwen2VL-7B and Qwen2VL-72B) as task proposers and autonomous evaluators, 371 thus completely eliminating the dependence of PAE on proprietary models. Results are included 372 in Table 1, 2, 3, and Figure 3. In particular, on WebArena, We found that LLava-7B PAE using 373 Qwen2VL-72B as the task proposer and evaluator achieved a similar performance as using Claude 3 374 Sonnet as the task proposer and evaluator, despite their significant difference in agent performances 375 (23.9% compared to 26.0% average success rate). As a result of this improvement, LLava-7B PAE using Qwen2VL-72B as the task proposer and evaluator achieved a better performance compared 376 to Qwen2VL-72B itself. Perhaps more surprisingly, even Qwen2VL-7B with much inferior agent 377 performance compared to LLaVa-7B SFT (7.5% compared to 18.0%) can be used to make significant improvements (from 18.0% to 23.1%). A similar conclusion is observed on WebVoyager
experiments as well, where even Qwen2VL (with agent performance of only 1.4%) can be used as
task proposers and evaluators to improve the performance of LLaVa7B-SFT from 14.9% on seen
websites and 9.1% on unseen websites to 21.7% and 13.7% on unseen websites. These results
demonstrate that the improvements from PAE root in the asymmetric capabilities of state-of-the-art
VLMs as agents and as task proposers/evaluators, instead of imitating a stronger VLM.

Does PAE scale well with larger and more capable base models? To test the scaling performance
 of PAE, we repeat our experiments on WebVoyager with a larger and more capable VLM base model
 LLaVa-1.6-34B (Liu et al., 2024). With a better base model, we still find a similar performance
 gain of PAE despite the model size change from 7B to 34B (7.4% compared to 10.8% absolute
 success rate improvement). Again, LLaVa-34B PAE beats LLaVa-7B PAE on 12 out of 13 websites
 from WebVoyager. Our scaling experiments suggest that PAE a favorable scaling property that can
 similarly improve better and larger base VLM agents as they become available.

391 Do the skills learnt by PAE generalize to unseen environments? To understand the generalization
of LLaVa-7B PAE to the websites that it has never interacted with, we apply the workflow from He
et al. (2024) to generate 500 tasks using Claude 3 Sonnet on 85 unseen online websites and test
the checkpoints from WebVoyager experiments. Results are presented in Table 3 and a list of the
websites is included in Appendix E. We observe that PAE for both LLaVa-7B and LLava-34B enable
the agents to learn general web-browsing skills that can be zero-shot transferred to unseen websites,
with 7.2% and 5.3% improvement in absolute success rate respectively.

6 DISCUSSIONS

398 399

400

401 The effect of additional reasoning step. We also perform an additional 402 ablation on the effect of the PAE de-403 sign choice of asking the VLMs to 404 output their thoughts first prior to the 405 actual web operations. We consider 406 an additional baseline of directly out-407 putting the web operations without 408 thoughts, and carry out the similar 409 SFT and Filtered BC experiments us-410 ing the same setup described in Sec-



Figure 3: Ablation experiments on WebArena Easy. The left figure measures the performance on the set of proposed tasks by different models with autonomous evaluators while the right figure measures the performance on WebArena Easy with the ground-truth evaluator.

tion 5.2. As reported in Figure 3, although PAE without reasoning can also achieve improvements
in the proposed set, the lack of additional reasoning step results in a significant inferior performance
in its generalization to the unseen human-written evaluation set.

414 The effect of choice of evaluators. Finally, we present ablation results on the effect of different 415 design choices of evaluators in Figure 3. We compare the outcome-based evaluator included in PAE 416 with other choices of evaluators in the related literature such as step-based evaluators (Pan et al., 417 2024) and function-based evaluators (Zhang et al., 2024a; Faldor et al., 2024). In our implementation of step-based evaluator, we ask Claude 3 Sonnet to evaluate whether each step is correct or not (i.e. 418 whether it gets the agent closer to the goal) and behavior clone all the steps considered correct by 419 the step-based evaluator. In our implementation of function-based evaluator, we provide 3 examples 420 of verification functions as used by WebArena (Zhou et al., 2024a) and ask Claude 3 Sonnet to also 421 come up with verification functions to functionally verify the final task success rate (e.g. checking 422 if the final website url is the same as the ground-truth url). As shown in Figure 3, both step-based 423 evaluator and function-based evaluator perform worse than the outcome-based evaluator, where the 424 use of step-based evaluator even leads to a worse performance compared to the SFT checkpoint 425 to start with. We found that the step-based evaluator hallucinated more often and tended to be 426 too "generous" in terms of considering the success of each step, potentially because the task of 427 evaluating the success of each step is significantly harder compared to only evaluating the success of 428 the final outcome. Furthermore, we found that oftentimes the function-based evaluator hallucinates on the success criterion for the verification function (e.g. making up a non-existing url that the 429 agent needs to go to), therefore resulting in most tasks being impossible to learn. In contrast, the 430 design choice of an outcome based evaluator can best provide reliable reward signals for the policy 431 to improve, resulting in better performances.

444



(a) Correlation between Human and Autonomous Evaluator. (b) Confusion Matrix

Figure 4: Correlation and confusion matrix analysis of different models in Webvoyager. (a) Correlation
 between human evaluations and our autonomous evaluator across various models at the system level. (b)
 Confusion matrix of the overall correlation between human evaluations and our autonomous evaluator at the
 instance level. Both results show strong correlation between our autonomous evaluator and human evaluations.

Alignment with human judgements. We demonstrate the effectiveness of our autonomous evaluator with a user study. We randomly select 200 trajectories for each method and present all screenshots in the trajectories, the corresponding actions at each step, and the task descriptions to the human annotator to decide if the task has actually been completed or not. As shown in Figure 4(a), there is a high correlation between our evaluator and human assessments across different models with an average misalignment of 1.7% at the system level and 8.9% at the instance level. The effectiveness of PAE as judged by human annotators is consistent with what is reported in Table 1.

- Error analysis. To understand where the improvement of PAE comes from, we conducted a user 452 study to analyze different error types across various models. With a high-level evaluation of model 453 capacities, we classified the error types into the following categories: (1) Low-level skills missing 454 **error** refer to the cases where the agent has a reasonable plan to solve the problem but fails to 455 execute precise actions on the website, such as not knowing which button to click to navigate to 456 the desired page. (2) High-level planning or reasoning errors refer to the cases where the agent 457 fails to generate a plan in its thoughts to solve the given task or cannot arrive at the correct answer 458 through reasoning with the website's screenshots. (3) Visual hallucinations refer to the cases where 459 the agent generates responses with made-up information that are not supported by the screenshot. 460 For example, the agent may claim that it has found a product that the task asked for while remaining at the homepage of google search, or the agent may produce a wrong answer while being on the right 461 page. (4) Timeouts refer to the cases where the agent is on the right track to solving the tasks, but 462 couldn't complete the task within maximum number of steps. (5) Technical Issues are not the fault 463 of the agent but caused by environment problems such as websites out of service and connection 464 issues. (6) Others include other less often error types such as the task itself is impossible. 465
- We present the results of error analysis for different models in WebVoyager in Figure 5, with a more 466 detailed analysis and full trajectories in Appendix I. Comparing LLaVa-7B SFT with LLaVa-34B 467 SFT, we observe that the predominant failure mode for LLaVa-7B SFT is visual hallucinations while 468 that for LLaVa-34B is low-level skill missing errors. This is because the reasoning capability for 469 LLaVa-7B base model is limited so it tends to imitate the demonstration data to produce answers 470 that look similar without being aware of the correctness of the answers. While LLaVa-34B SFT 471 is more aware of the correctness of answers (evidenced by a reduced visual hallucination rate), it 472 does not have the low-level web navigation skills so often fall short of low-level operations. PAE 473 can effectively improve on the major failure mode for both 7B and 34B models. In particular, for 474 LLaVa-7B SFT, PAE can reduce the visual hallucination rate (from 37% to 23%), making the agent 475 more aware of the goal of actually completing the tasks instead of imitating the demonstrations. 476 For LLaVa-34B SFT, PAE can effectively enrich the skill repertoire with low-level web navigation 477 skills, thereby reducing the low-level skill missing error (from 45% to 21%). Comparing our models with other VLM agents, we find that other open-source VLM agents such as Qwen2VL-7B and 478 Qwen2VL-72B mostly struggle with low-level web navigation skills while the error types for more 479 advanced proprietary models such as Claude 3.5 Sonnet are more spread out. 480

Comparison of different choices of contexts. We present our study on the effects of using different contexts on WebArena in Table 2 and Figure 6. By comparing the success rate between LLaVa-7B
 PAE and LLaVa-7B PAE (User Demos), we find additional information significantly improves the performance in the original WebArena task set Map (19.5% to 21.7%) but does not make a big difference in the rewritten easier task sets on PostMill and OneStopMarket. By manually inspecting the tasks proposed with and without user demos, we find that many tasks proposed with website



Figure 6: Online sample complexity comparisons on different websites in WebArena Easy between PAE using different contexts. Note that PAE with different contexts for task proposers uses different training tasks. Learning curves are smoothed with exponential running averages.

502 names alone are too hard or even impossible given the supported features of OpenStreetMap. For example, a task like "Locate the closest movie theater to the address 456 Oak Street, Chicago, Illi-504 nois, and provide the theater's name, address, and current movie showtimes." is impossible to be 505 completed on OpenStreetMap as it does not contain information related to the current movie showtimes. As shown in the learning curve in Figure 6, indeed the agent achieves a significantly lower 506 performance on the training tasks of PAE compared to that of PAE (User Demos). On the contrary, 507 this gap in terms of training set performances is much reduced on PostMill and OneStopMarket. We 508 hypothesize that this is because our simplified tasks on PostMill and OneStopMarket only examine 509 the basic usages of the websites such as "Go to a forum related to relationship advice" and "Browse 510 the Patio and Garden shopping category" and such tasks can be easily proposed with rudimentary 511 understanding of the websites inferred from the names of the websites alone. As the tasks get harder 512 and involve more complicated interactions with elements on different websites, we expect the use of 513 context information to play a more important role.

514 Qualitative comparisons. To qualitatively understand the benefits of PAE, we present snippets 515 of example trajectories in Figure 12 from evaluations on WebVoyager where LLaVa-7B PAE and 516 LLaVa-7B SFT attempt the same tasks. Full trajectories are included in Appendix I. In the first 517 example, we find that while LLaVa-7B knows SFT that it should use the search bar to find models 518 related to error correction, it fails to choose the correct search bar (should be [18] instead of [1]). 519 However, LLaVa-7B PAE learns the skill of using the search bar through typing into the correct 520 index [1] and executes its plan to complete the task. In the second example, the agent needs to 521 navigate to the Advanced Security page of Github. While both models are able to navigate to the 522 Security page of Github first, there turns out to be no direct links from the Security page to the Advanced Security page. As a result, LLaVa-7B SFT ends up wandering in Github without finding 523 the Advanced Security page. In contrast, LLaVa-7B PAE learns the skill of using Google Search in 524 the absence of a direct link and it successfully navigates to the right page with its help. In both cases, 525 we observe qualitative evidence of PAE teaching the agent a diverse repertoire that can effectively 526 help the agent to complete unseen tasks. 527

528 7 CONCLUSIONS AND FUTURE WORK

529 In this paper, we introduced a working system, PAE, for autonomous skill discovery with founda-530 tion model agents, addressing the limitations of using a static set of human-annotated instructions for fine-tuning agents. Instead of manually specifying what the agents should learn, our system 531 enables the agents to explore, practice, and refine new skills autonomously through open-ended 532 interactions with various environments. The framework's key components-task proposer, action 533 policy, and autonomous evaluator-work together to generate, attempt, and evaluate tasks without 534 any human intervention, leading to more than 10% improvement over prior state-of-the-art perfor-535 mance across benchmarks like WebArena Easy and WebVoyager among open-source VLM agents 536 (22.6% to 33%). This work paves the way for more capable open-source foundation model agents, with future research focused on extending this approach to other domains and integrating it with 538 better approaches to make use of the context information.

539

499

500

540 REPRODUCIBILITY STATEMENT

541 542 543

544

546

To facilitate reproducibility of our work, we plan to open-source the model checkpoint and code. To provide more details about our practical algorithm, we have included the algorithm pseudo-code in Algorithm 1. We have also included all the prompts that we have used for the task proposer, the agent policy, and the autonomous evaluator in Appendix C. More details for gathering and processing the SFT dataset have been included in Appendix E. An discussion of the hyperparameter tuning of our method has been included in Appendix H.

547 548 549

550

558

563

564

565

566 567

568

569

578

579

580 581

582

583

584

ETHICS STATEMENT

This work aims to enhance autonomous Internet agents through open-ended interactions with the web. However, the irresponsible or unrestricted use of such agents may pose risks, including personal data leaks or vulnerabilities to malicious attacks. To mitigate these risks, it is crucial to implement robust precautionary measures. In our experiments involving open-ended web navigation, we ensure that the agent is restricted from accessing personal accounts and employ appropriate firewalls to block DNS requests to suspicious websites. These safeguards help prevent unintended consequences and protect sensitive information.

- 559 REFERENCES
- Joshua Achiam, Harrison Edwards, Dario Amodei, and Pieter Abbeel. Variational option discovery algorithms, 2018. URL https://arxiv.org/abs/1807.10299.
 - Anthropic. The claude 3 model family: Opus, sonnet, haiku, 2024. URL https://www-cdn. anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_ Card_Claude_3.pdf.
 - Hao Bai, Yifei Zhou, Mert Cemri, Jiayi Pan, Alane Suhr, Sergey Levine, and Aviral Kumar. Digirl: Training in-the-wild device-control agents with autonomous reinforcement learning, 2024. URL https://arxiv.org/abs/2406.11896.
- 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.
 - Víctor Campos, Alexander Trott, Caiming Xiong, Richard Socher, Xavier Giro i Nieto, and Jordi Torres. Explore, discover and learn: Unsupervised discovery of state-covering skills, 2020. URL https://arxiv.org/abs/2002.03647.
 - 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.
- Cédric Colas, Tristan Karch, Nicolas Lair, Jean-Michel Dussoux, Clément Moulin-Frier, Peter Ford
 Dominey, and Pierre-Yves Oudeyer. Language as a cognitive tool to imagine goals in curiosity driven exploration, 2020. URL https://arxiv.org/abs/2002.09253.
- 588
 589
 589
 589
 590
 Cédric Colas, Laetitia Teodorescu, Pierre-Yves Oudeyer, Xingdi Yuan, and Marc-Alexandre Côté. Augmenting autotelic agents with large language models, 2023. URL https://arxiv.org/ abs/2305.12487.
- 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.

594 Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek 595 Gupta, and Jacob Andreas. Guiding pretraining in reinforcement learning with large language 596 models, 2023. URL https://arxiv.org/abs/2302.06692. 597 Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: 598 Learning skills without a reward function, 2018. URL https://arxiv.org/abs/1802. 06070. 600 601 Maxence Faldor, Jenny Zhang, Antoine Cully, and Jeff Clune. Omni-epic: Open-endedness via 602 models of human notions of interestingness with environments programmed in code, 2024. URL 603 https://arxiv.org/abs/2405.15568. 604 Hiroki Furuta, Kuang-Huei Lee, Ofir Nachum, Yutaka Matsuo, Aleksandra Faust, Shixiang Shane 605 Gu, and Izzeddin Gur. Multimodal web navigation with instruction-finetuned foundation models, 606 2024. URL https://arxiv.org/abs/2305.11854. 607 608 GeminiTeam. Gemini: A family of highly capable multimodal models, 2024. URL https:// 609 arxiv.org/abs/2312.11805. 610 Izzeddin Gur, Hiroki Furuta, Austin V. Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and 611 Aleksandra Faust. A real-world webagent with planning, long context understanding, and pro-612 gram synthesis, 2021. 613 614 Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan, 615 and Dong Yu. Webvoyager: Building an end-to-end web agent with large multimodal models, 616 2024. URL https://arxiv.org/abs/2401.13919. 617 Joey Hong, Sergey Levine, and Anca Dragan. Zero-shot goal-directed dialogue via rl on imagined 618 conversations, 2023a. URL https://arxiv.org/abs/2311.05584. 619 620 Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan 621 Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent: 622 A visual language model for gui agents, 2023b. URL https://arxiv.org/abs/2312. 623 08914. 624 Mengkang Hu, Pu Zhao, Can Xu, Qingfeng Sun, Jianguang Lou, Qingwei Lin, Ping Luo, Saravan 625 Rajmohan, and Dongmei Zhang. Agentgen: Enhancing planning abilities for large language 626 model based agent via environment and task generation, 2024. URL https://arxiv.org/ 627 abs/2408.00764. 628 629 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik 630 Narasimhan. Swe-bench: Can language models resolve real-world github issues?, 2024. URL 631 https://arxiv.org/abs/2310.06770. 632 Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham 633 Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating 634 multimodal agents on realistic visual web tasks, 2024. URL https://arxiv.org/abs/ 635 2401.13649. 636 637 Michael Laskin, Hao Liu, Xue Bin Peng, Denis Yarats, Aravind Rajeswaran, and Pieter 638 Abbeel. Unsupervised reinforcement learning with contrastive intrinsic control. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), Advances in Neu-639 ral Information Processing Systems, volume 35, pp. 34478-34491. Curran Associates, Inc., 640 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/ 641 file/debf482a7dbdc401f9052dbe15702837-Paper-Conference.pdf. 642 643 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023a. URL 644 https://arxiv.org/abs/2304.08485. 645 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 646 Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. URL https:// 647

llava-vl.github.io/blog/2024-01-30-llava-next/.

661

662

663

664

673

679

680

681

682

684

686

687

688

689

- 648 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, 649 Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui 650 Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. Agentbench: Evaluating llms as agents, 2023b. URL https://arxiv.org/abs/ 651 652 2308.03688.
- Llama3Team. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407. 654 21783. 655
- 656 OpenAI. Gpt-4 technical report, 2024. URL https://arxiv.org/abs/2303.08774. 657
- Jiavi Pan, Yichi Zhang, Nicholas Tomlin, Yifei Zhou, Sergey Levine, and Alane Suhr. Autonomous 658 evaluation and refinement of digital agents, 2024. URL https://arxiv.org/abs/2404. 659 06474. 660
 - Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho, He He, Sainbayar Sukhbaatar, and Jason Weston. Iterative reasoning preference optimization, 2024. URL https://arxiv.org/ abs/2404.19733.
- Seohong Park, Jongwook Choi, Jaekyeom Kim, Honglak Lee, and Gunhee Kim. Lipschitz-665 constrained unsupervised skill discovery, 2022. URL https://arxiv.org/abs/2202. 666 00914. 667
- 668 Seohong Park, Kimin Lee, Youngwoon Lee, and Pieter Abbeel. Controllability-aware unsupervised 669 skill discovery, 2023. URL https://arxiv.org/abs/2302.05103. 670
- Seohong Park, Oleh Rybkin, and Sergey Levine. Metra: Scalable unsupervised rl with metric-aware 671 abstraction, 2024. URL https://arxiv.org/abs/2310.08887. 672
- Pranav Putta, Edmund Mills, Naman Garg, Sumeet Motwani, Chelsea Finn, Divyansh Garg, and 674 Rafael Rafailov. Agent q: Advanced reasoning and learning for autonomous ai agents, 2024. 675 URL https://arxiv.org/abs/2408.07199. 676
- 677 Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware 678 unsupervised discovery of skills, 2020. URL https://arxiv.org/abs/1907.01657.
 - Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models, 2023a. URL https://arxiv.org/abs/2305.16291.
- 683 Tianlu Wang, Ilia Kulikov, Olga Golovneva, Ping Yu, Weizhe Yuan, Jane Dwivedi-Yu, Richard Yuanzhe Pang, Maryam Fazel-Zarandi, Jason Weston, and Xian Li. Self-taught eval-685 uators, 2024a. URL https://arxiv.org/abs/2408.02666.
 - Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. Cogvlm: Visual expert for pretrained language models, 2024b. URL https://arxiv.org/ abs/2311.03079.
- 691 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and 692 Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions, 2023b. URL https://arxiv.org/abs/2212.10560. 693
- 694 Tianhao Wu, Weizhe Yuan, Olga Golovneva, Jing Xu, Yuandong Tian, Jiantao Jiao, Jason Weston, 695 and Sainbayar Sukhbaatar. Meta-rewarding language models: Self-improving alignment with 696 llm-as-a-meta-judge, 2024. URL https://arxiv.org/abs/2407.19594. 697
- Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Sil-699 vio Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. Osworld: Benchmarking multimodal 700 agents for open-ended tasks in real computer environments, 2024. URL https://arxiv. 701 org/abs/2404.07972.

- 702 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, 703 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, 704 Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jin-705 gren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, 706 Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wen-707 bin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng 708 Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, 709 Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report, 2024a. URL 710 https://arxiv.org/abs/2407.10671. 711
- John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. Swe-agent: Agent-computer interfaces enable automated software engineering, 2024b. URL https://arxiv.org/abs/2405.15793.
- Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. Self-rewarding language models, 2024. URL https://arxiv.org/ abs/2401.10020.
- Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. Agenttun ing: 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, 2023a. URL https://arxiv.org/abs/2312.13771.
- Jenny Zhang, Joel Lehman, Kenneth Stanley, and Jeff Clune. Omni: Open-endedness via models of human notions of interestingness, 2024a. URL https://arxiv.org/abs/2306.01711.
- Jesse Zhang, Jiahui Zhang, Karl Pertsch, Ziyi Liu, Xiang Ren, Minsuk Chang, Shao-Hua Sun, and Joseph J. Lim. Bootstrap your own skills: Learning to solve new tasks with large language model guidance, 2023b. URL https://arxiv.org/abs/2310.10021.
- Pan Zhang, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Rui Qian, Lin Chen, Qipeng Guo, Haodong Duan, Bin Wang, Linke Ouyang, Songyang Zhang, Wenwei Zhang, Yining Li, Yang Gao, Peng Sun, Xinyue Zhang, Wei Li, Jingwen Li, Wenhai Wang, Hang Yan, Conghui He, Xingcheng Zhang, Kai Chen, Jifeng Dai, Yu Qiao, Dahua Lin, and Jiaqi Wang. Internlm-xcomposer-2.5: A versatile large vision language model supporting long-contextual input and output, 2024b. URL https://arxiv.org/abs/2407.03320.
- Zhuosheng Zhang and Aston Zhang. You only look at screens: Multimodal chain-of-action agents,
 2024. URL https://arxiv.org/abs/2309.11436.
 - Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. Gpt-4v(ision) is a generalist web agent, if grounded, 2024. URL https://arxiv.org/abs/2401.01614.
- 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, 2024a. URL https://arxiv.org/abs/
 2307.13854.
- Zhiyuan Zhou, Pranav Atreya, Abraham Lee, Homer Walke, Oier Mees, and Sergey Levine. Autonomous improvement of instruction following skills via foundation models, 2024b. URL https://arxiv.org/abs/2407.20635.

740

741

742

- 752
- 753

754

Appendices

756

757 758 759

760 761

762

A ALGORITHM

In Algorithm 1, we include a formal definitions of our practical algorithm of PAE as presented in Section 3.

764 Algorithm 1 Proposer-Agent-Evaluator: Practical Algorithm 765 **Require:** Context information $z_{\mathcal{M}}$, task proposer $\hat{\mathcal{C}}$, autonomous evaluator $\hat{\mathcal{R}}$. 766 1: Initialize policy π from a pre-trained checkpoint. 767 2: Initialize replay buffer $\mathcal{D} \leftarrow \{\}$. 768 3: ## Propose tasks based on the context information. 769 4: Obtain proposal task distribution $C(z_{\mathcal{M}})$. 770 5: for each global iteration do 771 6: for each trajectory to be collected do 772 7: Sample a task from the task proposer $c \sim \hat{C}(z_{\mathcal{M}})$. 773 Reset the environment to obtain the initial observation s_0 8: 774 9: for each environment step t do 775 10: Sample $a_t \sim \pi(\cdot | s_t, c), s_{t+1} \sim \mathcal{T}(\cdot | s_t, a_t, c).$ 776 11: if done then 777 ## Autonomously evaluate the outcome of the agent rollout. 12: 778 13: $r_t \leftarrow \mathcal{R}(s_t, a_t, c).$ 779 14: else $r_t \leftarrow 0.$ 15: end if 16: 781 17: $\mathcal{D} \leftarrow \mathcal{D} \cup \{(s_t, a_t, r_t, s_{t+1}, c)\}.$ 782 18: end for 783 19: end for 784 20: ## Update the agent policy with any RL algorithm. 785 $\pi \leftarrow \mathsf{RL}_update(\pi, \mathcal{D})$ 21: 786 22: end for 787

788

789

790 791

B DETAILS ON HUMAN ANNOTATIONS AND USER DEMOS

Details on User Demos. User demos from experiments in Figure 6 are collected by the authors without appealing to actual users. For each website, the authors attempt 10 tasks that we think are representative of the use of the particular website, and 10 most distinctive web pages are identified in the process of attempting those 10 tasks.

Details on Human Annotations. Five annotators of PhD students participate in the user study and the entire process takes around 40 annotator hours with the help of a designated user interface programmed in Gradio. To clarify the precise definition of the different error categories used in Section 6, we provide the following instruction to give more comprehensive explanations with example trajectories:

(1) Low-level skill missing errors refer to cases where the agent has a reasonable plan to solve the problem but fails to execute precise actions on the website, such as not knowing which button to click to reach the desired page. We classify trajectories where the agent seems to follow a reasonable plan but struggles with specific operations into this category. For example, in Figure 7, the task is "Find the Easy Vegetarian Spinach Lasagna recipe on Allrecipes and tell me what the latest review says." The agent attempts to search for the desired item but fails to click the correct button to reach the detailed page in the search results.

(2) High-level planning or reasoning errors occur when the agent fails to generate a complete plan
 or cannot reason correctly with the website's screenshots to solve the task. Trajectories where the agent cannot devise a plan for complex tasks or misinterprets the screenshot's content are categorized

as such. For instance, in Figure 8, the task is "Give 12 lbs of 4-cyanoindole, converted to molar and indicate the percentage of C, H, N." The agent should first search on Google about the chemical definition of 4-cyanoindole, then use WolframAlpha to calculate the result. However, the agent fails to get the precise definition of 4-cyanoindole, and doesn't know how to solve the task.

(3) Visual hallucinations refer to instances where the agent generates fabricated responses not supported by the screenshot. The agent might, for example, claim to have found a requested product while still on the Google homepage or provide an incorrect answer even when on the correct page. In Figure 9, the task is "Find out the trade-in value for an iPhone 13 Pro Max in good condition on the Apple website". The agent claims with a very detailed answer but actually it never access any page related to the trade-in on the website.

(4) Timeouts occur when the agent is on the right track to solving the task but cannot complete it within the maximum number of steps. This error indicates that the agent did nothing wrong but was constrained by the environment's step limits. For example, in Figure 10, the task is "Go to the Plus section of Cambridge Dictionary, find Image quizzes, and complete an easy quiz about Animals. Tell me your final score." The agent reaches the maximum time step limit (10) while attempting to finish the quiz.

(5) **Technical issues** are not caused by the agent but by environmental problems, such as websites being down or connection failures. In Figure 11, the ChromeDriver crashes after a valid operation.

(6) Others include less frequent error types, such as when the task itself is impossible to complete.

829 830 831

832

826

827

828

C ALL PROMPTS IN THE EXPERIMENTS

833 For completeness, we include examples of the prompts that we have used in this section. In partic-834 ular, in Figure 13, we have provided the prompt that we used for the Claude-Sonnet-3 autonomous 835 evaluator to evaluate the success for the task completion for all tasks in WebArena. A similar is used 836 for all tasks in WebVoyager. In Figure 14, 15, 16 we have included the prompts that we used for 837 generating the proposal tasks for each domain. We used the same prompts with 3 additional website 838 screenshots appended to the messages for PAE + User Demos. It is worth noting that our task pro-839 posers are domain-general and have little domain customizations. In particular, for all 13 real-world websites from WebVoyager, we use the same prompt to generate tasks except with the placeholder 840 of "web_name". This shows that our PAE framework can easily scale to multiple websites without 841 the need for domain-specific knowledge. The prompt for zero-shot VLM agents are included in 842 Figure 17, 18, and 19. 843

844 845

846 847

848

849

850

D PROMPTS FOR ZERO-SHOT VLM AGENTS

We also append the prompts (Figure 17, 18, and 19) that we used for the zero-shot baselines including Claude-Sonnet-3, Claude-Sonnet-3.5, Qwen2VL, InternVL2b5, LLaVa-1.6-7B, and LLaVa-1.6-34B. The prompt for WebVoyager tasks largely follow from that used in the prior literature (He et al., 2024). We include additional necessary domain knowledge of the WebArena tasks and evaluation protocols in the prompt that we used for WebArena.

851 852 853

854

E DETAILS FOR SFT

855 SFT for WebVoyager. As shown in Table 1, unlike proprietary VLMs, none of the open-source 856 VLM agent is able to follow the instructions and achieve non-trivial performances in real-world web navigation tasks in the zero-shot manner. Such models can rarely get success rewards in the process 858 of RL, thus leading to very slow convergence. To "warm-up" the open-source VLM agent to achieve 859 a non-trivial performance at the start of RL training, we turn to enhancing the performances with 860 SFT before RL. Note that the SFT process may not be needed if the base VLM agent model can 861 already achieve non-trivial performances such as Claude 3 Sonnet. To prevent data contamination, we gather 85 out-of-distribution real-world websites (listed in Figure 20 and 21), and collect 11220 862 trajectories in total using Claude 3 Sonnet with the prompt specified in Figure 17. The average 863 trajectory success rate is 25% as measured by our Claude 3 Sonnet evaluator. Each action in the



Figure 7: Extra full trajectories of fail trajectory 1, with error type Low-level Operational error, executed by model LLaVa-7B SFT. The task is 'Find the Easy Vegetarian Spinach Lasagna recipe on Allrecipes and tell me what the latest review says'. 916

914



Figure 8: Extra full trajectories of fail trajectory 2, with error type **Planning or Reasoning error**, executed by model LLaVa-7B PAE. The task is 'Give 12 lbs of 4-cyanoindole, converted to molar and indicate the percentage of C, H, N'.

trajectories contains both thoughts and actual web actions shown in Figure 2. All 11220 trajectories are used for SFT. RL training is carried out on top of the SFT checkpoint.

SFT for WebArena. In our preliminary experiments, we found that the SFT checkpoint trained on 946 real-world websites do not generalize well to self-hosted websites on WebArena. This is potentially 947 because of the distribution shift between real-world commercial websites and self-hosted websites. 948 For example, most real-world map websites such as Google Maps and Apple Maps support advanced 949 fuzzy search capabilities such as "pittsburgh to new york" while OpenStreetMap from WebArena 950 will not return any results with such queries. Therefore, we collect 3000 Claude 3 Sonnet generated 951 trajectories each from OpenStreetMap, Reddit, and OneStopMarket websites from WebArena. We 952 use the prompts from Figure 18 and 19 for the Claude agent. The average trajectory success rate is 953 27% as measured by our Claude 3 Sonnet evaluator. The SFT checkpoint for WebArena is fine-tuned 954 from the SFT checkpoint for WebVoyager.

955 956 957

938

939

944

945

F ADDITIONAL RESULTS ON WEBARENA

958 959 960

For completeness, we have also provided additional experiment results of different models from Ta-961 ble 2 in the original task split of WebArena (Zhou et al., 2024a). As shwon in the comparison results 962 presented in Table 4, even SOTA proprietary VLM agents like Claude 3 Sonnet struggle with the 963 tasks in WebArena with a success rate of only 14.6% with set-of-marks observations and chain-of-964 thought prompting. After performing SFT using the demonstrations generated by Claude 3 Sonnet, 965 LLaVa-7B SFT can only achieve 1.4% and 5.8% success rate on PostMill and OneStopMarket. By 966 manually inspecting the roll-out trajectories generated by LLaVa-SFT, we found that around half of 967 the successful trajectories on those two websites are false positives from the WebArena evaluator. 968 In these trajectories, the agent simply guessed the answer to be "no" or "N/A" where the ground truth happens to be that the task is not executable. As a result, the actual success rate on those two 969 websites is lower than 2%, leaving very sparse reward signals for RL to make meaningful improve-970 ments. We therefore rewrote the tasks on PostMill and OneStopMarket to be easier and report the 971 performances of PAE in Table 2.



Figure 9: Extra full trajectories of fail trajectory 3, with error type **Visual Hallucination**, executed by model LLaVa-7B SFT. The task is 'Find out the trade-in value for an iPhone 13 Pro Max in good condition on the Apple website'.

		OpenStreetMap	PostMill	OneStopMarket	Averag
Proprietary	Claude 3 Sonnet	24.3	10.6	11.2	14.6
	Qwen2VL-7B	0.7	0.0	1.3	0.7
Open-source	InternVL2.5-8B	2.6	0.2	3.3	2.3
•	LLaVa-7B	0.0	0.0	0.0	0.0
Ours	LLaVa-7B SFT	15.2	1.4	5.8	7.2

Table 4: Success rate comparisons across different domains from WebArena. Success and failure are detected with ground-truth verification functions. All tasks from OpenStreetMap are kept unchanged from WebArena task splits.



Figure 10: Extra full trajectories of fail trajectory 4, with error type **Timeouts**, executed by model Claude 3.5 Sonnet. The task is 'Go to the Plus section of Cambridge Dictionary, find Image quizzes and do an easy quiz about Animals and tell me your final score'.



Figure 11: Extra full trajectories of fail trajectory 5, with error type **Technical issues**, executed by model LLaVa-7B PAE. The task is 'Identify a course on Coursera that provides an introduction to Psychology, list the instructor's name, the institution offering it, and how many hours it will approximately take to complete'.



Figure 12: Qualitative comparison between LLaVa-7B PAE and LLaVa-7B SFT on the same tasks.
 LLaVa-7B PAE model successfully completed two tasks using learned skills from the RL training.

1157 1158 G LIMITATIONS

1159 Despite the progress of PAE for open-source VLM agents, there are still some limitations due to 1160 practical constraints. First of all, due to the limitations in fundamental capabilities of open-source 1161 base VLM models, our models trained with PAE are still inferior to state-of-the-art proprietary mod-1162 els in realistic web navigations, where advanced reasoning and planning capabilities are required. 1163 Moreover, because of the hallucination issues of open-source VLMs, we found them unreliable to 1164 serve as the autonomous evaluators and had to rely on advanced proprietary VLMs for judging the 1165 success and providing rewards. Finally, because of the dynamic nature of the real websites that we 1166 are using, some of our results may not be produced exactly, although a significant improvement from PAE should still be observed. 1167

1168

1169 H HYPERPARAMETERS

We include the hyperparameters that we have used in Table 5. As shown in the table, the only hyperparameters that PAE have on top of standard supervised fine-tuning are number of trajectories to collect in each global iteration in Algorithm 1, number of proposed tasks from the task proposer before RL training, and the number of seen screenshots for the evaluator. In our experiments, we found that PAE is relatively not sensitive to the choices of these hyperparameters, showing the robustness of PAE.

1177

¹¹⁷⁸ I MORE QUALITATIVE EXAMPLES

In this section, we present additional qualitative examples of agent trajectories while performing
 tasks to further demonstrate the effectiveness of our PAE. We will also release the full dataset for
 further analysis.

- Full trajectories of examples in Section 6. Here, we provide the complete trajectories for the examples discussed in the qualitative comparisons in Section 6, as shown in Figures 22–25. We detail the agent's thoughts and actions at each time step throughout the entire trajectory.
- **Some representative successful trajectories.** We also showcase representative successful trajectories generated by the LLaVa-7B PAE model to highlight the strengths of our method. In Figure 26,

1188 1189 1190 1191 1192 **Autonomous Evaluator Prompt** 1193 You are an expert in evaluating the performance of a web navigation agent. The agent is designed to 1194 help a human user navigate a website to complete a task. Your goal is to decide whether the agent's 1195 1196 execution is successful or not. 1197 As an evaluator, you will be presented with three primary components to assist you in your role: 1198 1. Web Task Instruction: This is a clear and specific directive provided in natural language, detailing 1199 the online activity to be carried out. 2. Result Response: This is a textual response obtained after the execution of the web task. It serves 1201 as textual result in response to the instruction. 1202 3. Result Screenshots: This is a visual representation of the screen showing the result or intermediate 1203 1204 state of performing a web task. It serves as visual proof of the actions taken in response to the 1205 instruction. 1206 - You SHOULD NOT make assumptions based on information not presented in the screenshot when 1207 comparing it to the instructions. 1208 - Your primary responsibility is to conduct a thorough assessment of the web task instruction against 1209 the outcome depicted in the screenshot and in the response, evaluating whether the actions taken 1210 1211 align with the given instructions. 1212 - NOTE that the instruction may involve more than one task, for example, locating the garage and 1213 summarizing the review. Failing to complete either task, such as not providing a summary, should 1214 be considered unsuccessful. 1215 - NOTE that the screenshot is authentic, but the response provided by LLM is generated at the end 1216 of web browsing, and there may be discrepancies between the text and the screenshots. 1217 - NOTE that if the content in the Result response is not mentioned on or different from the screen-1218 shot, mark it as not success. 1219 1220 You should explicil consider the following criterions: 1221 - Whether the claims in the response can be verified by the screenshot. E.g. if the response claims the 1222 distance between two places, the screenshot should show the direction. YOU SHOULD EXPECT 1223 THAT THERE IS A HIGH CHANCE THAT THE AGENT WILL MAKE UP AN ANSWER NOT 1224 VERIFIED BY THE SCREENSHOT. 1225 - Whether the agent completes EXACTLY what the task asks for. E.g. if the task asks to find a 1226 specific place, the agent should not find a similar place. 1227 1228 In your responses: You should first provide thoughts EXPLICITLY VERIFY ALL THREE CRITE-1229 RIONS and then provide a definitive verdict on whether the task has been successfully accomplished, 1230 either as 'SUCCESS' or 'NOT SUCCESS'. 1231 A task is 'SUCCESS' only when all of the criteria are met. If any of the criteria are not met, the task 1232 should be considered 'NOT SUCCESS'. 1233 1234 Figure 13: The prompt used by the autonomous evaluator for Claude-Sonnet-3. Same prompt is 1235 used to evaluate tasks from WebArena websites. The evaluator takes as inputs the task description, 1236 the response from the agent's ANSWER action, and last three screenshots in the trajectory. The 1237 evaluation result is a binary verdict of 'SUCCESS' or 'NOT SUCCESS'. 1238 1239

1240

1242	Task Proposer Prompt for WebVoyager
1243	{"web_name": "Apple", "id": "Apple-40", "ques": "Find the pricing and specifications
1244	for the latest Mac Studio model, including the available CPU and GPU options.", "web":
1246	"https://www.apple.com/"}
1247	We are training a model to navigate the web. We need your help to generate instructions. With the
1248	examples provided above, please give 25 more example tasks for the model to learn from in the
1249	domain of {web_name}. You should imagine tasks that are likely proposed by a most likely user of
1250	this website. A few demos of users navigating through the web are provided above.
1251	YOU SHOULD MAKE USE OF THE DEMOS PROVIDES TO GENERATE TASKS, SO THAT
1252	YOUR TASKS ARE REALISTIC AND RELEVANT TO THE WEBSITE.
1253	Please follow the corresponding guidelines:
1255	1)First output your thoughts first on how you should come up with diverse tasks that examine various
1256	capabilities on the particular website, and how these tasks reflect the need of the potential user. Then
1257	you should say 'Output:' and then followed by the outputs STRUCTURED IN JSONL FORMAT.
1258	You should not say anything else in the response.
1259	2)PLEASE MAKE SURE TO HAVE 25 examples in the response!!!
1260	3)Your proposed tasks should be DIVERSE AND COVER A WIDE RANGE OF DIFFERENT
1262	POSSIBILITIES AND DIFFICULTY in the domain of {web_name}. Remember, your job is to
1263	propose tasks that will help the model learn to navigate the web to deal with various real world
1264	requests.
1265	4)Your task should be objective and unambiguous. The carry-out of the task should NOT BE DE-
1266	PENDENT on the user's personal information such as the CURRENT TIME OR LOCATION.
1267	5)You should express your tasks in as diverse expressions as possible to help the model learn to
1260	understand different ways of expressing the same task.
1270	6)Your tasks should be able to be evaluated OBJECTIVELY. That is, by looking at the last three
1271	screenshots and the answer provided by an agent, it should be possible to tell without ambiguity
1272	whether the task was completed successfully or not.
1273	7)Your tasks should require a minimum completion steps from 3 to 7 steps, your tasks should have a
1274	diverse coverage in difficulty as measured by the minimum completion step. I.E. You should propose
1275	not only tasks that may take more than 4 steps to complete but also tasks that can be completed within
1270	3 steps.
1278	8)Humans should have a 100% success rate in completing the task.
1279	9)Your tasks should be able to be completed without having to sign in to the website.
1280	
1281	Figure 14: Prompts used by Claude-Sonnet-3 for proposing tasks in WebVoyager experiments. For
1282	PAE + User Demos, we use the same prompt with additional user demos appended to the message.
1203	
1285	the task is "Show the most played games on Steam, and tell me the number of players currently
1286	in-game." In Figure 27, the task is "Find out the starting price for the most recent model of the iMac
1287	on the Apple website. In Figure 28, the task is 'Look up the use of modal verbs in the grammar section for expressing possibility (e.g. 'might' 'could' 'may') and find examples of their usage
1288	in sentences on the Cambridge Dictionary." Finally, in Figure 29, the task is "Search for plumbers"
1289	available now but not open 24 hours in Orlando, FL."
1290	
1292	
1293	

	Task Proposer Prompt for WebArena Map
	{"web_name": "map", "id": "map-2", "ques": "Tell me the full address of all international ai
	hat are within a driving distance of 50 km to University of California, Berkeley"}
	{"web_name": "map", "id": "map-10", "ques": "I will arrive San Francisco Airport soon. Pr
	he name of a Hilton hotel in the vicinity, if available. Then, tell me the the shortest walking di
	to a supermarket from the hotel."}
	{"web_name": "map", "id": "map-17", "ques": "Check if the ikea in pittsburgh can be reach
	one hour by car from hobart street"}
	We are training a model to navigate the web. We need your help to generate instructions. We
	examples provided above, please give 25 more example tasks for the model to learn from
	domain of OpenStreetMap. You should imagine who is the most likely user for the websi
	propose tasks that are likely to be proposed by this user. Please follow the corresponding guide
	1)First output your thoughts first on how you should come up with diverse tasks that examine v
	capabilities on the particular website, and how these tasks reflect the need of the potential user.
	you should say 'Output:' and then followed by the outputs STRUCTURED IN JSONL FOR
	You should not say anything else in the response.
	2)PLEASE MAKE SURE TO HAVE 25 examples in the response!!!
	3)Your proposed tasks should be DIVERSE AND COVER A WIDE RANGE OF DIFFE
	POSSIBILITIES AND DIFFICULTY in the domain of OpenStreetMap. Remember, your jo
	propose tasks that will help the model learn to navigate the web to deal with various real wo
(quests. 4)Your task should be objective and unambiguous. The carry-out of the task should NC
	DEPENDENT on the user's personal information such as the CURRENT TIME OR LOCAT
	5)You should express your tasks in as diverse expressions as possible to help the model le
	understand different ways of expressing the same task.
	5)Your tasks should be able to be evaluated OBJECTIVELY. That is, by looking at the last
	screenshots and the answer provided by an agent, it should be possible to tell without amb
	whether the task was completed successfully or not.
	7)Your tasks should require a minimum completion steps from 3 to 7 steps, your tasks should
	diverse coverage in difficulty as measured by the minimum completion step. I.E. You should pr
	not only tasks that may take more than 4 steps to complete but also tasks that can be completed y
	not only tasks that may take more than 4 steps to complete but also tasks that can be completed s 3 steps.
	not only tasks that may take more than 4 steps to complete but also tasks that can be completed 3 steps. 3 steps. 3)Humans should have a 100% success rate in completing the task.

To de Deren e and Deren Welt America De dall'America Ora Claus Marchest
Task Proposer Prompt for WebArena Reddit and OneStopMarket
{"web_name": "Apple", "id": "Apple=40", "ques": "Find the pricing and specifications
for the latest Mac Studio model, including the available CPU and GPU options.", "web":
"https://www.apple.com/"}
We are training a model to navigate the web. We need your help to generate instructions. With the
examples provided above, please give 25 more example tasks for the model to learn from in the
domain of {web_name}.
You should provide tasks in the DOMAIN OF {web_name}.
Please follow the corresponding guidelines: 1)First answer how many screenshots are provided and
describe in detail the functions of the website that you see from each of the screenshot. Then output
your thoughts first. Then you should say 'Output:' and then followed by the outputs STRUCTURED
IN ISONI, FORMAT. You should not say anything else in the response
2)PLEASE MAKE SURE TO HAVE 25 examples in the response!!!
4)Your tack should start from the home page of the website instead of the shown screenshots
4) Four task should start from the nome page of the website instead of the should even in diverse.
5) Four task does not need to be the same as real users would do, but it should examine diverse
capabilities of the agent to do web navigartion.
6)Your tasks should examine the VERY BASIC functions of the website and should not require
complicated web page operations. They can be completed within 5 steps.
7)THIS DOMAIN IS A SELF-HOSTED STATIC DOMAIN AND DIFFERENT FROM POPULAR
WEBSITES, DO NOT ASSUME ANY INFORMATION NOT PROVIDED IN THE SCREEN-
SHOTS.
8)Your tasks should examine the capability of the web agent to find some information on the website,
navigating to some specific web pages. Do not propose tasks that involve making actual modifica-
tions to the websites.
9)Your tasks should result in the agent landing in a single groundtruth web page or finding a single
grounth truth answer. The landed webpage can be some specific categories, a drafted post, some
search results, or even the homepage of the website. When the task is to to find some information.
specify exactly what information the agent should find such as the price, the number of comments.
the title etc. It can also be information about the current account
Figure 16: Prompts used by Claude-Sonnet-3 for proposing WebArena Tasks for Reddit and On-
eStopMarket. For PAE + User Demos, we use the same prompt with additional user demos ap-
pended to the message.

Zero-Shot VLM Agent Prompt for Web Voyager $(1/2)$ Imagine you are a robot browsing the web just like humans. Now you need to complete a task. In
each iteration, you will receive an Observation that includes a screenshot of a webpage and some
texts. This screenshot will feature Numerical Labels placed in the TOP LEFT corner of each Web
to the Web Element that requires interaction then follow the guidelines and choose one of the
following actions:
1. Click a Web Element.
3. Scroll up or down. Multiple scrolls are allowed to browse the webpage. Pay attention!! The
default scroll is the whole window. If the scroll widget is located in a certain area of the webpage,
then you have to specify a Web Element in that area. I would have the mouse there and then scroll.
5. Go back, returning to the previous webpage.
6. Google, directly jump to the Google search page. When you can't find information in some
websites, try starting over with Google. 7 Answer This action should only be chosen when all questions in the task have been solved
7. This wei. This action should only be chosen when an questions in the disk have been solved.
Correspondingly, Action should STRICTLY follow the format:
- Chek [Numerical_Label] - Type [Numerical_Label]; [Content]
- Scroll [Numerical_Label or WINDOW]; [up or down]
- Wait
- Goble
- ANSWER; [content]
Kay Chidalines You MUST fallow
* Action guidelines *
1) To input text, NO need to click textbox first, directly type content. After typing, the system
automatically hits 'ENTER' key. Sometimes you should click the search button to apply search
2) You must Distinguish between textbox and search button, don't type content into the button! If
no textbox is found, you may need to click the search button first before the textbox is displayed.
 3) Execute only one action per iteration. 4) STRICTLY Avoid repeating the same action if the webpage remains unchanged. You may
have selected the wrong web element or numerical label. Continuous use of the Wait is also NOT
allowed.
5) When a complex Task involves multiple questions or steps, select "ANSWER" only at the very end after addressing all of these questions (steps). Flexibly combine your own abilities with
the information in the web page. Double check the formatting requirements in the task when
ANSWER.
1) Don't interact with useless web elements like Login. Sign-in. donation that appear in Webpages
Pay attention to Key Web Elements like search textbox and menu.
2) Vsit video websites like YouTube is allowed BUT you can't play videos. Clicking to download
3) Focus on the numerical labels in the TOP LEFT corner of each rectangle (element). Ensure vou
don't mix them up with other numbers (e.g. Calendar) on the page.
4) Focus on the date in task, you must look for results that match the date. It may be necessary to find the correct year month and day at calendar
5) Pay attention to the filter and sort functions on the page, which, combined with scroll, can help
you solve conditions like 'highest', 'cheapest', 'lowest', 'earliest', etc. Try your best to find the
answer that best fits the task.
Your reply should strictly follow the format:
Thought: {Your brief thoughts (briefly summarize the info that will help ANSWER)}
Action: {One Action format you choose}
Then the User will provide:
Observation: {A labeled screenshot Given by User}

1455 LLaVa-1.6-34B.

1458	
1459	
1460	
1461	Zero-Shot VLM Agent Prompt for WebArena (2/2)
1462	each iteration, you will receive an Observation that includes a screenshot of a webnage some texts
1463	and the accessibility tree of the webpage. This screenshot will feature Numerical Labels placed in
1464	the TOP LEFT corner of each Web Element. The accessibility tree contains information about the
1465	web elements and their properties. The numrical labels in the screenshot correspond to the web
1/66	elements in the accessibility tree.
1/67	Carefully analyze the visual information to identify the Numerical Label corresponding to the
1407	Web Element that requires interaction, then follow the guidelines and choose one of the following
1408	actions:
1469	2 Delete existing content in a textbox and then type content
1470	3. Scroll up or down. Multiple scrolls are allowed to browse the webpage. Pay attention!! The
1471	default scroll is the whole window. If the scroll widget is located in a certain area of the webpage.
1472	then you have to specify a Web Element in that area. I would have the mouse there and then scroll.
1473	4. Wait. Typically used to wait for unfinished webpage processes, with a duration of 5 seconds.
1474	5. Go back, returning to the previous webpage.
1475	6. Answer. This action should only be chosen when all questions in the task have been solved.
1476	Correspondingly Action should STRICTLY follow the format:
1477	- Click [Numerical Label]
1478	- Type [Numerical Label]: [Content]
1479	- Scroll [Numerical_Label or WINDOW]; [up or down]
1480	- Wait
1481	- GoBack
1482	- ANSWER; [content]
1483	Kay Guidalinas You MUST follow:
1484	* Action guidelines *
1/05	1) To input text, NO need to click textbox first, directly type content. After typing, the system
1405	automatically hits 'ENTER' key. Sometimes you should click the search button to apply search
1400	filters. Try to use simple language when searching.
1487	2) You must Distinguish between textbox and search button, don't type content into the button! If
1488	no textbox is found, you may need to click the search button first before the textbox is displayed.
1489	3) Execute only one action per iteration.
1490	4) STRICTLY Avoid repeating the same action if the wedpage remains unchanged. You may have selected the wrong web element or numerical label. Continuous use of the Wait is also NOT
1491	allowed
1492	5) When a complex Task involves multiple questions or steps, select "ANSWER" only at the
1493	very end, after addressing all of these questions (steps). Flexibly combine your own abilities with
1494	the information in the web page. Double check the formatting requirements in the task when
1495	ANSWER.
1496	6) If you can't find the answer using the given website because there is no such information on the
1497	website after some attempts, you should report "N/A" as the answer to represent that the task is impossible to solve with the given webgite. You may have 15 store to try to solve the task.
1498	7) Only provide answer based on the information from the image make sure the answer is consistent
1499	with the image, don't hallucinate any information that is not based on image.
1500	
1501	* Web Browsing Guidelines *
1502	1) Focus on the numerical labels in the TOP LEFT corner of each rectangle (element). Ensure you
1503	don't mix them up with other numbers (e.g. Calendar) on the page.
150/	2) Pay attention to the filter and sort functions on the page, which, combined with scroll, can help
1505	you solve conditions like highest, cheapest, lowest, earliest, etc. Try your best to find the
1505	
1000	
1507	

Figure 18: The prompt used for all zero-shot VLM agents for WebArena websites, including Claude-Sonnet-3, Claude-Sonnet-3.4, Qwen2-VL, InternVL-2.5-XComposer, LLaVa-1.6-7B, and LLaVa-1.6-34B. To be continued in Figure 19.

	Zero-Shot VLM Agent Prompt for WebArena
	* OpenStreetMap Usage Guidelines *
	1) When you need to search the address of a location, you can just type the location in the 'search'
	bar. You don't need to use the directions button to get the address. The directions button is only
	used when you need to find the distance/walk/drive time between two locations.
	OpenStreetMap does not support approximate search. You may get no results. This is because the
	keywords or try to find the location by yourself. Note that openstreet map does not support search
	phrase like 'Cafe near CMU", you should try to find it by yourself.
	3) When you need to find the distance/walk/drive time between two locations, you should FIRST
	CLICK ON THE DIRECTIONS BUTTON (drawn as two arrows), to the right of the 'Go' Button and usually labeled as [10] or [11] AND ONLY INDUTTING THE TWO LOCATIONS AFTER
	CLICKING ON THE DIRECTIONS BUTTON WHEN THE DIRECTIONS SEARCH BARS
	ARE SHOWN.
	4) When you are trying to type some locations in the directions search bar, sometimes you may
	receive an alert of 'couldn't locate' followed by the location you typed. This means the location
	you typed is not found in the map. Do not immediately try something else. You need to quit the
	5) When you search the walk/drive/bike time, make sure that you are USING THE RIGHT MODE
	OF TRANSPORTATION. The default mode is usually set to 'Drive'.
	6) When you need to get the DD of some location, you need to click the location shown in the
	search result in the left part of the screen. The DD will be shown then starting with 'Location:'.
	7) When you need to answer the zip code of some location, you should directly answer the 5-digit
	Zip code. The answer should be " 15232 " instead of "The Zip code of the location is 15232 ". Note that the zincode will be displayed in the search result, you don't need to click the location to the
	information page to find the zip code.
	8) When you need to answer the phone of some location, please omit the part of the country code.
	The answer should be "4122683259" instead of "+1 412 268 3259".
	* Reddit Usage Guidelines *
	screenshot. You do not need to further payigate to the reddit website
	2) When you want to find a subreddit, you need to first navigate to Forums to see the list of
	subreddits. Under forums, you will see only a subset of subreddits. To get the full list of subreddits,
	you need to navigate to the Alphabetical option. To know you can see the full list of subreddits, you
	will see 'All Forums' in the observation. Often you will not find a focused subreddit that exactly
	have reached a subreddit successfully you will see 'f/subreddit name' in the observation
	3) When you want to post forum in reddit, remember to fill up all the content, then click the button
	'Create forum'. The button maybe located below out of the screenshot, you need to scroll down to
	find it.
	4) When you want to ask or post something in a subreddit, you need to first find that subreddit and
	then finish the work. 5) forums and subreddits are the same thing.
	Your reply should strictly follow the format:
	Thought: Your brief thoughts (briefly summarize the info that will help ANSWER)
	Action: One Action format you choose
	רי ווי דד ו, וייי
	I nen the User will provide:
	Remember only execute one action in each step. For example 'Action: Type [8]: CMU Type [9]
	Pittsburgh' is not allowed. You should execute the action 'Type [8]; CMU' first. then 'Type [9]
	Pittsburgh' in the next step.
	Remember to always make your answer simple and clear. For example, if you want to report the zip
	code of some location, always say "AINSWER; 06516" instead of "The zip code of the location is 06516"
	00510 .
Fim	re 19: The prompt used for all zero-shot VLM agents for WebArena websites including Cl
וטויד	and the second second second second second for the second se

Out-of	listribution Websites of WebVoyager for SFT (1/2)
Allreci	es:
Simply	<pre>Recipes: https://www.simplyrecipes.com</pre>
Food N	twork: https://www.foodnetwork.com
Taste of	Home: https://www.tasteofhome.com
Yumml	: https://www.yummly.com
Food.co	n: https://www.food.com
Amazo	:
eBay: 1	tps://www.ebay.com
Walma	:https://www.walmart.com
Target:	<pre>ittps://www.target.com</pre>
Best Bi	/: https://www.bestbuy.com
Alibaba	https://www.alibaba.com
Apple:	
Samsur	: https://www.samsung.com
Micros	ft:https://www.microsoft.com
Sony: 1	tps://www.sony.com
Google	Store: https://store.google.com
Dell: h	tps://www.dell.com
ArXiv:	
SSRN:	ttps://www.ssrn.com
Researc	Gate: https://www.researchgate.net
bioRxiv	https://www.biorxiv.org
IEEE X	Nore: https://ieeexplore.ieee.org
Publie	nttps://pubmed.ncbi.nim.nin.gov
GitHul	
GitLab	https://about.gitlab.com
Bitbuck	t: https://bitbucket.org
Sourcel	orge: https://sourceforge.net
Codeba	e:https://www.codebasehq.com
Gitea: 1	ttps://gitea.io
ESPN∙	
CBS Sr	nts: https://www.cbssports.com
Fox Sn	rts: https://www.foxsports.com
NBC S	orts: https://www.nbcsports.com
Bleach	Report: https://www.bleacherreport.com
Sky Sp	rts: https://www.skysports.com
• 1	
Course	a:
edX: ht	tps://www.edx.org
Udacity	https://www.udacity.com
Udemy	https://www.udemy.com
Futurel	arn: https://www.futurelearn.com
⊾nan A	auemy: nttps://www.knanacademy.org

Figure 20: A list of 85 websites that we used to collect demonstration trajectories with Claude 3 Sonnet. In total 11220 trajectories were collected with different tasks. These websites were also used for testing the zeroshot generalization of PAE to out-of-distribution websites in Section 5. List continued in Figure 21.

6	out-of-distribution Websites of WebVoyager for SFT (2/2)
	ambridge Dictionary:
Ň	ferriam-Webster: https://www.merriam-webster.com
D	victionary.com: https://www.dictionary.com
С	xford Learner's Dictionaries: https://www.oxfordlearnersdictionaries.co
C	ollins English Dictionary: https://www.collinsdictionary.com
Y	ourDictionary: https://www.yourdictionary.com
R	BC Nows
C	NN: https://www.cnn.com
Ă	l Jazeera: https://www.aljazeera.com
R	euters: https://www.reuters.com
Т	he Guardian: https://www.theguardian.com
N	BC News: https://www.nbcnews.com
Ģ	oogle Maps:
A	pple Maps: https://maps.apple.com
N B	Ing Maps: https://www.bing.com/maps
IV V	Apquest. https://www.mapquest.com
v H	lere WeGo: https://wego_here_com
G	boogle Search:
B	ing: https://www.bing.com
Y	ahoo Search: https://search.yahoo.com
D	uckDuckGo: https://duckduckgo.com
B	aidu: https://www.baidu.com
Y	andex: https://yandex.com
H	lugging Face.
C	DenAI: https://openai.com
Ť	ensorFlow: https://www.tensorflow.org
Р	yTorch: https://pytorch.org
K	aggle: https://www.kaggle.com
S	paCy: https://spacy.io
v	Valfrom Alabor
v	vonram Aipna;
N	fathway: https://www.mathway.com
S	vmbolab https://www.svmbolab.com
N	ficrosoft Math Solver: https://mathsolver.microsoft.com
D	esmos: https://www.desmos.com



Figure 22: Full trajectories of success trajectory 1 in Figure 12 with task 'Find the most recently
updated machine learning model on Huggingface which focuses on Error Correction' executed by
model LLaVa-7B PAE.



Figure 23: Full trajectories of fail trajectory 1 in Figure 12 with task 'Find the most recently updated machine learning model on Huggingface which focuses on Error Correction' executed by model LLaVa-7B SFT.



Figure 24: Full trajectories of success trajectory 2 in Figure 12 with task 'Find the Security topic in GitHub Resources and answer the role of GitHub Advanced Security' executed by model LLaVa-7B PAE.



Figure 25: Full trajectories of fail trajectory 2 in Figure 12 with task 'Find the Security topic in GitHub Resources and answer the role of GitHub Advanced Security' executed by model LLaVa-7B SFT.

Figure 26: Extra full trajectories of successful trajectory 1 with task 'Show most played games in Steam. And tell me the number of players in In game at this time' executed by model LLaVa-7B PAE.

Figure 27: Extra full trajectories of successful trajectory 2 with task 'Find out the starting price for the most recent model of the iMac on the Apple website' executed by model LLaVa-7B PAE.

Figure 28: Extra full trajectories of successful trajectory 3 with task 'Look up the use of modal verbs in grammar section for expressing possibility (e.g., 'might', 'could', 'may') and find examples of their usage in sentences on the Cambridge Dictionary' executed by model LLaVa-7B PAE.

Figure 29: Extra full trajectories of successful trajectory 4 with task 'Search for plumbers available now but not open 24 hours in Orlando, FL' executed by model LLaVa-7B PAE.