

LEARNING TO CONTEXTUALIZE WEB PAGES FOR ENHANCED DECISION MAKING BY LLM AGENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent advances in large language models (LLMs) have led to a growing interest in developing LLM-based agents for automating web tasks. However, these agents often struggle with even simple tasks on real-world websites due to their limited capability to understand and process complex web page structures. In this work, we introduce LCoW, a framework for Learning language models to Contextualize complex Web pages into a more comprehensible form, thereby enhancing decision making by LLM agents. LCoW decouples web page understanding from decision making by training a separate contextualization module to transform complex web pages into comprehensible format, which are then utilized by the decision-making agent. We demonstrate that our contextualization module effectively integrates with LLM agents of various scales to significantly enhance their decision-making capabilities in web automation tasks. Notably, LCoW improves the success rates of closed-source LLMs (e.g., Gemini-1.5-flash, GPT-4o, Claude-3.5-Sonnet) by an average of 20%, and demonstrates a 33.5% average improvement in success rates for open-source LMs (e.g., Llama-3.1-8B, Llama-3.1-70B) on the WorkArena benchmark. Moreover, the Gemini-1.5-flash agent with LCoW achieves state-of-the-art results on the WebShop benchmark, outperforming human experts.

1 INTRODUCTION

Large language models (LLMs) have demonstrated strong potential in automating web tasks by treating web browsing as a sequential decision-making process, where web pages serve as observations and user interactions, such as clicking and typing, function as actions (Yao et al., 2022a;b). Various approaches have been developed to enhance the performance of LLM agents in these tasks. One such method involves fine-tuning open-source LLMs using demonstration data from web browsing tasks (Furuta et al., 2023; Lai et al., 2024). While promising, this approach requires extensive data collection and significant computational resources for effective fine-tuning. Alternatively, several studies have utilized advanced closed-source LLMs, such as GPT-4o (OpenAI, 2024), with carefully designed prompting techniques (Drouin et al., 2024; Zhou et al., 2023a; Sodhi et al., 2024; Pan et al., 2024). By leveraging the general world knowledge and reasoning capabilities of these models, the methods enhance web automation but at the cost of reduced controllability.

However, despite the advancements, state-of-the-art LLM agents often struggle to process complex raw web content, such as HTML and accessibility trees, posing significant challenges for their effective use in web task automation. While LLMs excel in tasks that require detailed reasoning, such as solving mathematical problems or coding, we hypothesize that their underperformance in seemingly

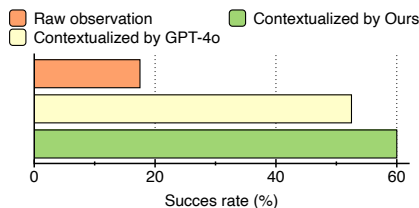


Figure 1: Success rate of the Gemini-1.5-flash agent on 40 WorkArena tasks. We selected a subset of 40 tasks by simply choosing the first 40 tasks based on the task indices. When the agent leverages observations contextualized by GPT-4o (yellow), its success rate improves by 31%, with further improvements achieved with our method (green).

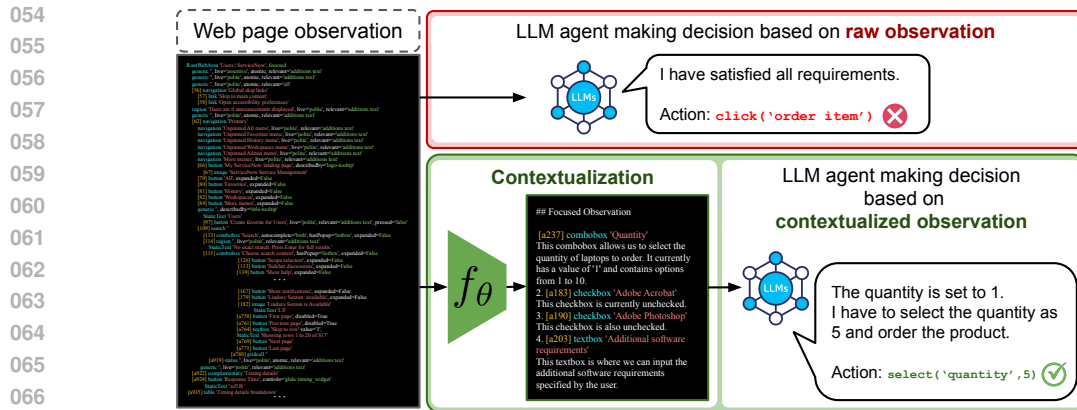


Figure 2: **(Top)** In the conventional pipeline, LLM agents decide on the next action based on raw, complex web page observations (e.g., HTML, accessibility trees), which often hinder accurate decision making. **(Bottom)** In our proposed pipeline, a contextualization module transforms these complex web page observations into a more comprehensible format, thereby enabling LLM agents to make more accurate decisions by enhancing their understanding of the web page.

simple decision-making tasks like web browsing is not due to a lack of decision-making capabilities but rather to difficulties in understanding and processing complex web page observations.

To validate our hypothesis, we conducted an initial experiment that demonstrated an LLM agent based on Gemini-1.5-flash can achieve substantial improvements in web browsing tasks when equipped with a module designed to contextualize complex web page observations (i.e., contextualization module). This module enhances task performance by removing irrelevant UI elements and highlighting key components with explanations, thereby simplifying the decision-making process. We evaluated the performance of this agent on 40 tasks from WorkArena (Drouin et al., 2024), a benchmark designed to assess web agents on real-world, enterprise-related websites. As shown in Figure 1, utilizing GPT-4o as the contextualization module (yellow) resulted in a 31% absolute improvement in the agent’s success rate compared to direct processing of raw observations (red). These results support our hypothesis that the difficulty in understanding web pages is a major bottleneck for LLM-based web agents.

In this work, we propose LCoW, a framework that includes a contextualization module and a training algorithm to fine-tune this module to enhance the decision-making capabilities of LLM agents in web automation. As illustrated in Figure 2, the contextualization module transforms complex web page observations into a comprehensible format, enabling LLM agents to make more accurate decisions. Furthermore, to enable the contextualization module to provide context more grounded in real websites, we propose an iterative algorithm designed to train the contextualization module. Our training algorithm consists of three phases: (i) trajectory collection, (ii) sampling contextualized observations, and (iii) updating the contextualization module. Notably, the proposed method does not depend on manually curated data for training; instead, it gathers data through the agent’s interactions within the web browsing environment. For each observation from the collected trajectories, we generate multiple contextualized observations using the current contextualization module. Each observation is then assigned a reward based on whether a set of LLM agents can accurately predict the correct action given the contextualized observation. Finally, we select the one with the highest reward as the target and train the contextualization module to maximize the likelihood of the target given the original raw observation.

As demonstrated in our initial experiment using the Gemini-1.5-flash (see Figure 1), LCoW significantly enhances the decision-making capabilities of LLM agents, even beyond the improvements seen with state-of-the-art LLMs like GPT-4o used as a contextualization module. In our experiments, we conduct comprehensive evaluations of our proposed approach on WebShop (Yao et al., 2022a) and WorkArena (Drouin et al., 2024), two popular benchmarks for evaluating agent performance in web environments. First, we demonstrate that LCoW significantly enhances the overall performance of LLM agents with varying scales (Llama-3.1-8B, Llama-3.1-70B, Gemini-1.5-flash, GPT-4o, Claude-3.5-Sonnet) across both benchmarks. In particular, LCoW achieves state-of-the-art

108 results on the WebShop benchmark, outperforming human experts. Second, we analyze how the
 109 contextualization module refines complex web pages and how the contextualization enhances the
 110 decision-making of LLM agents.

112 2 BACKGROUND

113 In this section, we describe the formulation of web browsing as a sequential decision-making prob-
 114 lem and the use of LLMs as decision-making agents.

115 Web browsing can be formulated as a Partially Observable Markov Decision Process (POMDP), de-
 116 fined by $\langle \mathcal{S}, \mathcal{O}, \mathcal{A}, T, \mathcal{R} \rangle$. The state $s_t \in \mathcal{S}$ represents the internal configuration of the web browser
 117 at time step t , which is only partially observable. The observation $o_t \in \mathcal{O}$ corresponds to the web
 118 page rendered by the web browser given s_t , which can take various forms (e.g., screenshot, HTML,
 119 accessibility trees). The action space \mathcal{A} is the set of all possible interactions with the UI elements
 120 (e.g., clicking, typing). The state transition function T defines the probability of transitioning from
 121 state s to state s' after performing an action $a_t \in \mathcal{A}$, such as clicking a link or scrolling. While tran-
 122 sitions are typically deterministic, occasional stochastic events (e.g., pop-ups, network errors) can
 123 occur. The reward function \mathcal{R} assesses the functional correctness, evaluating whether the resulting
 124 state s_t satisfies pre-defined criteria for successful task completion.

125 Leveraging their ability to interpret web pages and generate actions in text form, LLMs are increas-
 126 ingly employed as agents for automating web-based tasks. In this framework, an LLM agent π
 127 generates an action a_t to interact with UI elements at each time step t , based on the user instruc-
 128 tion [TASK], the current web page observation o_t , and the history of previous actions $a_{<t}$. The
 129 objective of the agent is to complete the given task in order to maximize the reward.

132 3 METHOD

133 In this section, we present LCoW, a framework for enhancing the capability of LLM agents by con-
 134 textualizing complex web pages. Section 3.1 outlines the concept of the contextualization module
 135 and its integration with LLM agents for web automation tasks. Section 3.2 introduces an iterative
 136 algorithm for training the contextualization modules to improve decision making of LLM agents.

137 3.1 CONTEXTUALIZATION MODULE

138 In this work, we decouple web page understanding from decision making of LLM agents. Our
 139 hypothesis is that while LLM agents possess strong decision-making capabilities, their performance
 140 can significantly degrade when the observations they rely on are lengthy and non-contextualized,
 141 such as HTML and accessibility trees. To address this limitation, we introduce a *contextualization*
 142 *module*, a separate language model designed to enhance LLM agents by contextualizing complex
 143 web page observations into a form that is more easily processed and comprehensible. Intuitively,
 144 a proper contextualization of observations can enhance the agent’s understanding of web content
 145 and its decision making based on that understanding (see Figure 3 for an example of the input and
 146 output of the module). Formally, given a web page observation o_t at time step t , the objective
 147 of our contextualization module f_θ is to generate a contextualized observation o_t^{co} that serves as
 148 input to the LLM agent π to enhance its decision making. Specifically, f_θ uses the task instruction
 149 [TASK], the previous actions of the agent $a_{<t}$, and the current web page observation o_t to generate
 150 a contextualized observation:

$$151 \quad 152 \quad 153 \quad 154 \quad o_t^{\text{co}} = f_\theta([\text{TASK}], a_{<t}, o_t).$$

155 The LLM agent π then predicts the next action based on the contextualized observation:

$$156 \quad 157 \quad 158 \quad 159 \quad a_t = \pi([\text{TASK}], a_{<t}, o_t^{\text{co}}).$$

160 While an arbitrary language model can serve as a contextualization module f_θ , it is important for
 161 the module to learn from experience in the web environment to provide more grounded context for
 decision making, such as role of a particular button or interaction with specific UI elements.

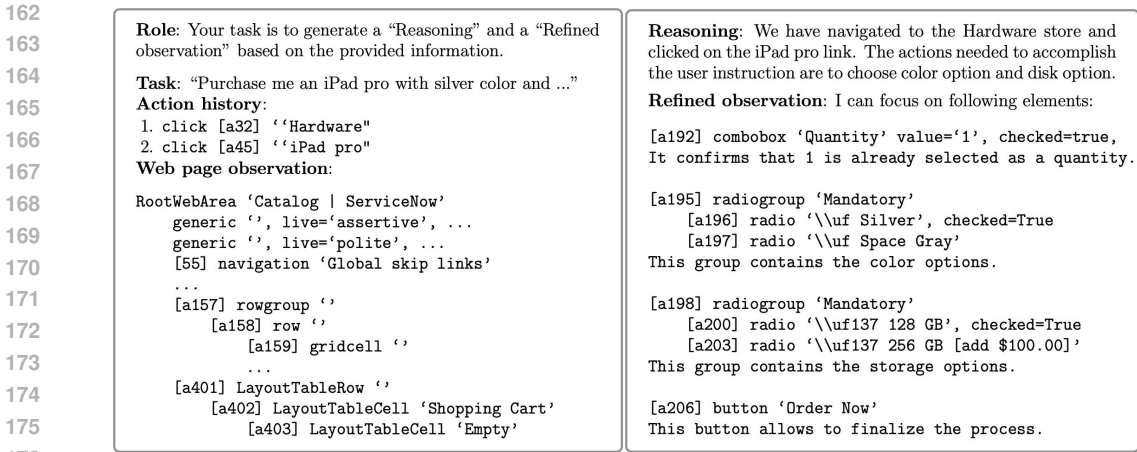


Figure 3: An example of a input of contextualization module including lengthy web page observation (left) and an observation contextualized by the contextualization module trained using LCoW (right). The module converts raw observations into a more concise form to enhance decision making in agents. The prompt used is provided in Appendix A.4.

3.2 ALGORITHM FOR TRAINING THE CONTEXTUALIZATION MODULE

We now describe an iterative algorithm for training the contextualization module f_θ to enhance decision making of LLM agents. In a nutshell, the algorithm involves an iterative process of collecting paired input-output data for training the contextualization module and subsequently updating the module based on the collected data. For data collection, we begin by gathering trajectories of successfully completed tasks from the web browsing environment. For each observation o_t in the collected trajectories, we sample multiple candidate contextualized observations, and select the one that best provides the relevant context for multiple LLM agents to accurately predict the next action a_t . Based on the chosen target observations, we update f_θ via supervised fine-tuning. We now outline a single iteration of LCoW, followed by a detailed explanation of the design of the reward used for evaluating the candidate contextualized observations.

Single iteration A single iteration of LCoW starts with the contextualization module $f_{\theta^{(i)}}$ and aims to update this module to $f_{\theta^{(i+1)}}$. This process consists of three phases:

Step 1 (Trajectory collection). Given a set of training tasks, we roll out the LLM agent π in the web environment to collect trajectory data. Specifically, the agent determines the next action based on the contextualized observation produced by $f_{\theta^{(i)}}$ until the episode terminates. We collect only those trajectories that end in the successful completion of the tasks.

Step 2 (Sampling optimal contextualization). As illustrated in Figure 4, we start by sampling multiple candidates o_t° from the current contextualization module $f_{\theta^{(i)}}$ (i.e., $o_t^{\circ} \sim f_{\theta^{(i)}}(\cdot \mid [\text{TASK}], a_{<t}, o_t)$) for each web page observation o_t in the collected trajectories. Each candidate is then assigned a reward based on whether a set of LLM agents can accurately predict the ground-truth action a_t given o_t , with the candidate receiving the maximum reward selected as the optimal contextualized observation. If all candidates receive a zero reward, we retry the sampling process with the ground-truth action a_t provided as additional context to $f_{\theta^{(i)}}$ (i.e., $o_t^{\circ} \sim f_{\theta^{(i)}}(\cdot \mid [\text{TASK}], a_{<t}, o_t, a_t)$) to guide the generation of valid contextualized observations.

Step 3 (Model update). We update the current module $f_{\theta^{(i)}}$ by fine-tuning it with the optimal contextualized observations collected in Step 2. With the updated module $f_{\theta^{(i+1)}}$, we return to Step 1 and repeat the process.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

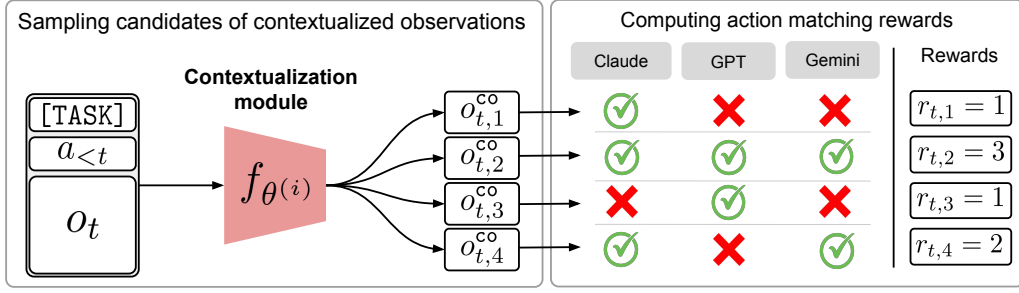


Figure 4: Illustration of sampling optimal contextualization. First, we sample multiple candidates of contextualized observations, given user instruction [TASK], previous actions $a_{<t}$, and observation o_t . Subsequently, multiple LLM agents predict the next action based on each candidate, and the reward for each candidate is computed according to how many LLM agents correctly predict the ground-truth action a_t . In the figure, $o_{t,2}^{co}$ receives the highest reward and is therefore used as the target data for updating $f_{\theta^{(i)}}$.

Algorithm 1 One iteration of LCoW

Require: a contextualization module $f_{\theta^{(i)}}$, an LLM agent π , a set of LLM agents for computing the reward $\Pi = \{\pi_i\}_{i=1}^K$, an empty trajectory buffer \mathcal{T} , an empty data buffer \mathcal{D} , and a set of training tasks \mathcal{G}_{tr}

- 1: // Trajectory collection
- 2: **for** [TASK] $\in \mathcal{G}_{tr}$ **do** ▷ Collect successful trajectories from training environment
- 3: $\tau, R \sim \pi(\cdot \mid [\text{TASK}], f_{\theta^{(i)}})$
- 4: **if** $R = 1.0$ **then**
- 5: $\mathcal{T}.\text{append}(\tau)$
- 6: // Sampling optimal contextualizations
- 7: **for** (o_t, a_t) in \mathcal{T} **do**
- 8: **for** $n \leftarrow 1$ to N **do** ▷ Sample N candidate contextualized observations and assign rewards
- 9: $o_{t,n}^{co} \sim f_{\theta^{(i)}}(\cdot \mid [\text{TASK}], a_{<t}, o_t)$
- 10: $r_{t,n} = \sum_{\pi \in \Pi} \text{ActionMatchingScore}(\pi([\text{TASK}], a_{<t}, o_{t,n}^{co}), a_t)$
- 11: **if** $\max_n(r_{t,n}) = 0$ **then**
- 12: **for** $n \leftarrow 1$ to N **do** ▷ Retry the sampling if rewards are zero for all candidates
- 13: $o_{t,n}^{co} \sim f_{\theta^{(i)}}(\cdot \mid [\text{TASK}], a_{<t}, o_t, a_t)$
- 14: $r_{t,n} = \sum_{\pi \in \Pi} \text{ActionMatchingScore}(\pi([\text{TASK}], a_{<t}, o_{t,n}^{co}), a_t)$
- 15: $o_{t,*}^{co} = \arg \max_n(r_{t,n})$
- 16: $\mathcal{D}.\text{append}([\text{TASK}], a_{<t}, o_t, o_{t,*}^{co})$
- 17: // Parameter update
- 18: $\theta^{(i+1)} := \arg \max_{\theta^{(i)}} \mathbb{E}_{([\text{TASK}], o_{t,*}^{co}, a_{<t}, o_t) \sim \mathcal{D}} [f_{\theta^{(i)}}(o_{t,*}^{co} \mid [\text{TASK}], a_{<t}, o_t)]$
- 19: **return** $f_{\theta^{(i+1)}}$

Algorithm 1 outlines the steps for a single iteration of the training algorithm. Starting with an initial contextualization module $f_{\theta^{(0)}}$, we iteratively train the module through M iterations. After each iteration, the module updates from $f_{\theta^{(i)}}$ to $f_{\theta^{(i+1)}}$ until reaching the final $f_{\theta^{(M)}}$.¹

Reward for contextualized observations The reward for the contextualized observation o_t^{co} is defined as the sum of the action-matching scores computed using multiple LLM agents. Each action matching score evaluates whether an LLM agent π correctly predicts the ground-truth action a_t given the contextualized observation o_t^{co} . By leveraging multiple LLM agents to compute the reward, we ensure that the module produces contextualized observations that generalize across a diverse set of agents, preventing overfitting to the behavior of any single agent and enhancing the module’s

¹In the initial iteration, we assume a limited set of seed demonstrations and use relatively strong LLMs to sample candidate contextualized observations to accelerate training.

270 adaptability to arbitrary LLM agents. Section 4.2 presents empirical results on the generalization
271 capabilities of the contextualization module.
272

273 4 EXPERIMENTS 274

275 We design our experiments to investigate the following questions:
276

- 277 • How effective is LCoW in training a contextualization module for improving decision making of
278 LLM agents?
- 279 • Can the contextualization module trained with LCoW generalize to arbitrary LLMs with varying
280 scales?
- 281 • What form do the web page observations take after contextualization and how the contextualized
282 observation aid the decision making of LLM agents?
283

284 4.1 EXPERIMENTAL SETUP 285

286 **Benchmarks** WebShop (Yao et al., 2022a) and WorkArena (Drouin et al., 2024) are two popular
287 benchmarks designed to evaluate the capabilities of web agents in completing various web tasks.
288 WebShop provides a simulated online shopping environment with real-world product data, consist-
289 ing of 500 evaluation tasks and 5,500 training tasks, each defined by natural language instructions
290 to purchase products that meet specific criteria. During each decision-making step, the agent re-
291 ceives HTML input to predict the next action, and at the end of an episode, it receives a reward
292 ranging from 0 to 1 based on how well the attributes of the purchased item align with the intended
293 product criteria. WorkArena, on the other hand, focuses on assessing web agents’ ability to complete
294 enterprise-related tasks, such as creating user accounts and ordering products from a service catalog,
295 on real-world websites. This benchmark comprises 33 task types, with each type containing up to
296 1,000 individual task instances, where agents must navigate websites by following natural language
297 instructions and receive rewards based on successful task completion.
298

299 **Training details** For the WebShop benchmark, we fine-tune Phi-3-mini-Instruct as a contextual-
300 ization module, setting the learning rate, warmup ratio, and batch size to $1e-5$, $1e-2$, and 32, respec-
301 tively. The module is trained for a single epoch over the collected data for each LCoW iteration, and
302 we utilize demonstrations corresponding to 397 individual tasks provided in the WebShop bench-
303 mark as seed demonstrations. In the WebShop environment, we employ Gemini-1.5-flash as the
304 LLM agent, combined with the contextualization module during the trajectory collection phase. It is
305 also utilized for sampling optimal contextualization at the initial iteration of LCoW. For WorkArena,
306 we fine-tune Llama-3.1-8B-Instruct as an observation contextualization module. Here, we set the
307 learning rate, warm-up ratio, and batch size as $1e-5$, $1e-1$, and 128 respectively, training the model
308 for 4 epochs. Since no demonstrations are provided in this benchmark, we collect 264 seed demon-
309 strations across 23 task types using Claude-3.5-Sonnet and GPT-4o, while no demonstrations were
310 collected for the remaining 10 task types. Summary statistics of the collected seed demonstrations
311 can be found in Appendix 8. Additionally, as determination of action matching based on parsing
312 is infeasible due to open-ended actions (e.g., sending message to user), we exploit GPT-4o as an
313 action-matching evaluator. Detailed prompt for the action-matching evaluator is provided in Ap-
314 pendix A.4. For WorkArena, we use Claude-3.5-Sonnet as the agent LLM for trajectory collection
315 and we utilize the same LLM for the sampling optimal contextualization at the initial iteration.

316 **Evaluation setup** In WebShop, we train the contextualization module using LCoW on environ-
317 ments associated with 500 of the 5,500 training tasks in WebShop. After training, we evaluate the
318 module on 500 evaluation tasks, measuring both the success rate and average reward, with a task
319 considered successful if the reward equals 1. To prompt the agent, we utilize a one-shot exemplar
320 as described in Yao et al. (2022b). The maximum number of state transitions is limited to 10, and
321 episodes are terminated if the same action is repeated more than twice. In the WorkArena, we sam-
322 ple 5 task instances for evaluation and 15 for training from each task type. As we were not able
323 to collect valid seed demonstrations from 10 among 33 task types, we conduct evaluation in 5 task
instances for each remaining 23 task types (i.e., 115 individual evaluation tasks as a total). Here,
the maximum number of state transitions is set to 20, and episodes are also terminated if the same

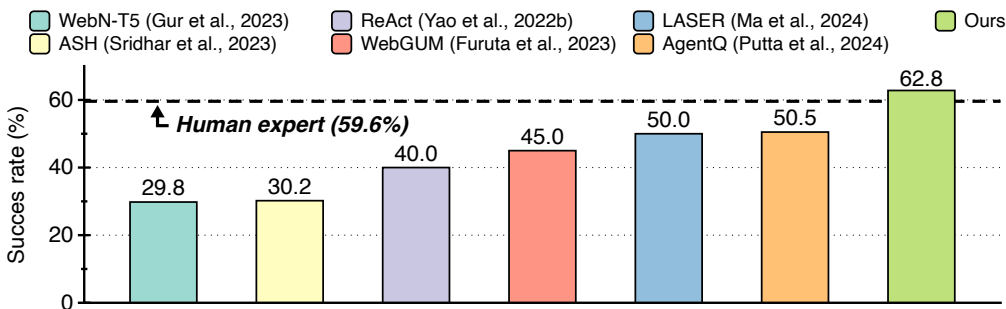


Figure 5: Success rate on 500 evaluation tasks from WebShop. Average human performance and expert human performance are 50% and 59.6%, respectively (Yao et al., 2022a). The Gemini-1.5-flash agent with the contextualization module trained for three iterations achieves a state-of-the-art success rate of 62.8%, outperforming the human expert performance, as well as previous baselines (Yao et al., 2022b; Furuta et al., 2023; Putta et al., 2024; Sridhar et al., 2023; Ma et al., 2024; Gur et al., 2023).

	GPT-4o		Gemini-1.5-flash		Claude-3.5-Sonnet		Llama-3.1-70B (Unseen)	
	Success	Reward	Success	Reward	Success	Reward	Success	Reward
Base prompt	34.8%	0.496	43.6%	0.693	26.6%	0.336	34.2%	0.590
Self-ctx	26.2%	0.459	46.4%	0.608	12.4%	0.146	40.2%	0.547
LCoW (iter 1)	27.8%	0.545	46.4%	0.705	39.4%	0.600	39.2%	0.666
LCoW (iter 2)	46.0%	0.647	58.2%	0.796	58.8%	0.780	55.0%	0.781
LCoW (iter 3)	50.6%	0.666	62.8%	0.803	59.8%	0.771	59.6%	0.803

Table 1: We investigate the efficacy of LCoW across multiple LLM agents in WebShop. For all LLM agents, LCoW consistently improves both success rate and reward over iterations, surpassing self-contextualization (self-ctx) and even human expert-level success rate by the third iteration. Additionally, LCoW is also effective when combined with Llama-3.1-70B, which was not used for computing the action-matching reward (i.e., unseen) during training the contextualization module.

action is repeated more than twice. For prompting the LLM agent, we use the prompt provided by BrowserGym (Drouin et al., 2024), a unified framework for the development and evaluation of web agents. The complete prompts used can be found in Appendix A.4.

4.2 MAIN RESULTS

We first demonstrate the effectiveness of LCoW on the WebShop and WorkArena benchmarks. Furthermore, we show that the contextualization module trained via LCoW enhances performance even for an LLM agent that was not involved in the LCoW training process, such as Llama-3.1-70B and Llama-3.1-8B. We evaluate against two baselines: (i) decision making based solely on raw observations without the contextualization module (i.e., base prompt), and (ii) decision making based on the observation contextualized by the LLM agent itself (i.e., self-contextualization).

Effectiveness of LCoW As shown in Table 1, both the GPT-4o and Claude-3.5-Sonnet agents perform poorly on the WebShop benchmark, frequently struggling to complete the tasks. Instead of selecting an item from the search results and reviewing its details, the agents often browse through multiple options in an attempt to find an exact match to the instructions, which frequently results in unsuccessful episode terminations. Although self-contextualization improves the success rate with Gemini-1.5-flash and Llama-3.1-70B, the performance of both the Claude-3.5-Sonnet and GPT-4o agents declines. In contrast, both the success rate and the average reward achieved by all three LLM agents improve substantially when integrated with the contextualization module trained with LCoW. Particularly, both the Gemini-1.5-flash and Claude-3.5-Sonnet agents surpass the average human performance of 50.5% when combined with the contextualization module trained for 2 iterations.

	GPT-4o	Gemini-1.5-flash	Claude-3.5-Sonnet	Llama-3.1-70B (Unseen)	Llama-3.1-8B (Unseen)
Base prompt	50.9%	14.2%	69.6%	36.8%	2.6%
Self-ctx	33.0%	17.4%	73.0%	29.6%	9.5%
LCoW (iter 1)	58.3%	59.2%	77.4%	56.5%	50.4%

Table 2: We evaluate the success rate of five LLM agents with varying scales on 115 tasks in the WorkArena benchmark, which is based on real-world website.

When the contextualization module is trained for 3 iterations, the agents exceed the expert human-level performance of 59.6%. As illustrated in Figure 5, the Gemini-1.5-flash agent combined with LCoW achieves state-of-art performance on the WebShop benchmark, outperforming prior methods in success rate by more than 12%. Notably, the Claude-3.5-Sonnet agent, lower performing than the other agents on WebShop, achieves superhuman performance at 59.8% when integrated with our contextualization module, demonstrating that LCoW can effectively enhance the decision-making capabilities of LLM agents.

As shown in Table 2, on WorkArena, the Claude-3.5-Sonnet agent outperforms the GPT-4o agent, achieving a success rate of 69.6% in our evaluation setup. When integrated with LCoW, Claude-3.5-Sonnet agent achieves an even higher success rate of 77.4%, which is higher than all baselines evaluated. GPT-4o and Gemini-1.5-flash also improve 7.4% and 45%, respectively, when combined with LCoW.

Generalization to arbitrary LLM agents We evaluate the Llama-3.1-8B agent and Llama-3.1-70B agent integrated with LCoW to assess its effectiveness with LLM agents not involved in the training process (i.e., those not used for computing action-matching rewards). As shown in Table 1, the Llama-3.1-70B agent combined with the contextualization module trained for two iterations outperforms average human performance (50.0%) on WebShop and achieves the success rate of human experts (59.6%) when trained for three iterations. Furthermore, Table 2 shows that the contextualization module also enhances Llama-3.1-8B and Llama-3.1-70B agent on WorkArena, improving the success rate by approximately 47% and 20%, respectively. It is noteworthy that a relatively small LLM agent (i.e., Llama 3.1-8B) struggles to perform tasks when given raw observations, but when combined with LCoW, its success rate rises on par to GPT-4o. This result further supports our hypothesis that the bottleneck in LLM agents performing web automation tasks lies in their observation understanding capability rather than decision-making capability. It also demonstrates that LCoW can elicit a significant level of decision-making capability even from smaller models at the 8B scale.

4.3 ANALYSIS

Can LCoW optimize the contextualization module for specific LLM agents? We demonstrate that LCoW effectively optimizes the contextualization module to generate context that leads to more accurate decision making by LLM agents. In the WebShop benchmark, we used the contextualization module checkpoint from each LCoW training iteration to generate contextualized observations for the 1,372 raw observations present in the 397 seed demonstrations. We then calculated the action-matching reward by comparing the actions predicted from these contextualized observations with the ground-truth actions derived from the demonstrations. Figure 6 shows the average action-matching reward across three iterations, demonstrating that with each round of training, the contextualization module learns to increasingly produce contextualized observations that enhance decision making of the LLM agents used in data collection (i.e., GPT-4o, Gemini-1.5-flash, Claude-3.5-Sonnet). This shows the effectiveness of LCoW in optimizing the contextualization module for specific LLM

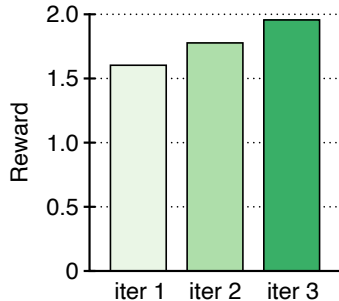


Figure 6: Average action-matching rewards across three iterations on WebShop show a consistent increase, suggesting that the contextualization module is optimized to generate observations that enhance decision-making in LLM agents.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

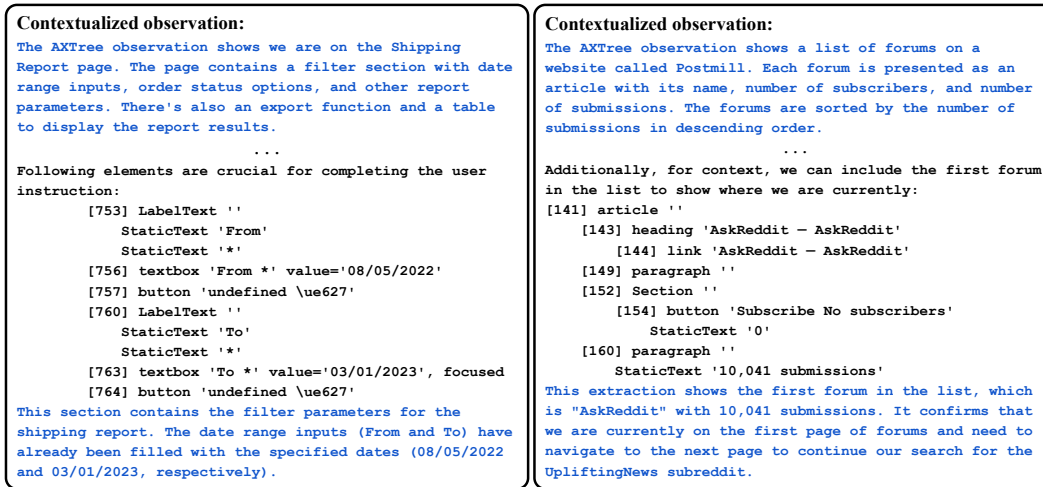


Figure 7: Examples of how the contextualization module refines complicated web pages into comprehensible format. As indicated by blue color, the contextualization module provides comprehensible context by verbally explaining the web page and UI elements relevant to the given task.

agents. Furthermore, it suggests the potential to extend LCoW as a method for indirectly tuning the behavior of closed-source LLM agents, where direct optimization is not feasible, across a wide range of text-based decision-making tasks.

How web pages are contextualized? We qualitatively analyze how the contextualization module, trained via LCoW, processes real-world web page observations and how the contextualization enhance decision making of LLM agents on web tasks. We observe that the contextualization module simplifies web page content by extracting UI elements relevant to the given task and providing clear descriptions of their functionalities and interaction methods. These contextualized explanations aid LLM agents significantly in making more accurate decisions, while reducing the selection of inadmissible or redundant actions, which is the main failure mode of the LLM agent for web browsing. For example, as shown in Figure 7, the contextualization module provides the explanation that to interact with a checkbox element, the LLM agent must click the associated `LabelText` element, an action that may not be immediately intuitive. Based on this information, the agent avoids redundant actions, such as repeatedly clicking the checkbox itself and instead makes more informed decisions. Additional examples are provided in the Appendix A.5. Such knowledge about the functionality and interaction of UI elements ground to the real website appears to have been learned during the LCoW training process, where the contextualization module explores various observation contextualizations and learns to generate contextualized observations that lead to more accurate decision making by LLM agents. Figure 8 shows that, with the contextualization module, the average number of decision-making steps by the agent decreases in both Claude-3.5-Sonnet and Llama-3.1-70B agent, allowing tasks to be completed more efficiently. More detailed examples are also provided in the Appendix A.5.

Comparison to directly training LLM agents One might argue that, given available demonstrations, it would be more straightforward to train the LLM agent directly rather than using the demonstrations to train a contextualization module. To demonstrate the superiority of training the contextualization module with an equivalent amount of demonstration data, we conduct comparative experiments against directly training LLM agent. In this analysis, we fine-tune Llama-3.1-8B with 264 seed demonstrations using behavior cloning (BC) as the baseline. To ensure a fair comparison in terms of model scale, we define a Llama-3.1-8B agent that performs tasks based on the output of a contextualization module trained using LCoW as the direct comparison group. As shown in Figure 9, both smaller LLM agents, such as Llama-3.1-8B, and larger LLM agents achieve higher success rates when using the LCoW-trained contextualization module compared to the BC baseline.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

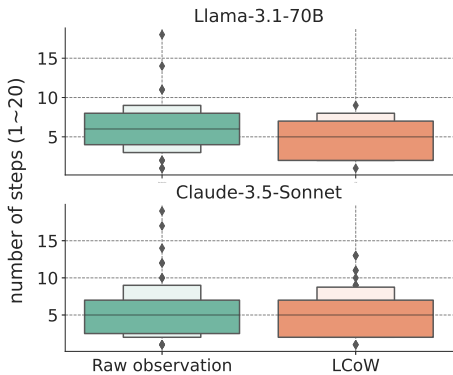


Figure 8: Comparison of the number of steps required to complete tasks in WorkArena. We only consider tasks that LCoW and baseline both succeed. Claude-3.5-Sonnet and Llama-3.1-70B agent both demonstrate efficient decision making by avoiding erroneous or redundant actions when the contextualization module is utilized.

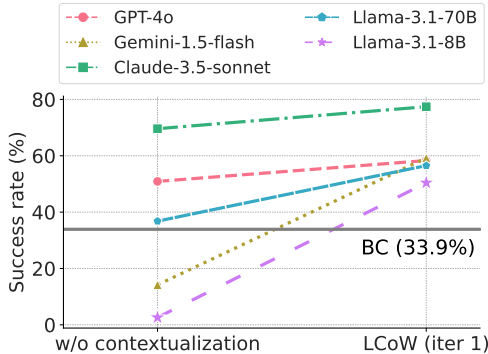


Figure 9: Comparison to training LLM agent via behavior cloning (BC). We fine-tune Llama-3.1-8B with demonstrations as a BC baseline. LLM agents combined with LCoW achieves much higher performance compared to BC baseline.

Limitations While the contextualization module trained using LCoW has shown to enhance LLM agents by converting web pages into a comprehensible format, it was hard to expect generalization to UI elements too specific to certain functionalities unseen during the training. More detailed analysis is provided in Appendix A.2. Although our experiment utilized small scale self-generated dataset for training the contextualization module, we believe that scaling up LCoW can be an interesting future direction. Additionally, we believe that trajectory collection can be further improved in order to collect successful trajectories from diverse open-ended tasks, thereby facilitating iterative training process. For example, applying search algorithms to collect more successful trajectories from complicated tasks or synthesizing novel tasks over iteration might be an interesting research direction.

5 RELATED WORK

LLM agents for web automation As LLMs continue to improve and demonstrate remarkable performance across diverse domains, developing web automation agents based on LLMs is gaining growing interest (Zhou et al., 2023a; Drouin et al., 2024). Recent works have explored various methods to utilize LLMs as agents for automating real-world web tasks. Pan et al. (2024) propose to apply a self-refine mechanism (Madaan et al., 2023) to improve decision making of agents through self-generated feedback. Sodhi et al. (2024) enhance web automation by using LLMs to manage low-level workflows handcrafted by humans, while Wang et al. (2024) extend this approach by introducing agent workflow memory, which extracts reusable routines from past experiences.

Similar to our work, a promising direction for enhancing web automation involves integrating a summarization module that condenses web page observations, enabling agents to predict actions based on these summarized inputs. For example, Deng et al. (2024) propose MindAct, which consists of an HTML extraction module and an action prediction module. The HTML extraction module ranks individual HTML elements using a ranking language model trained to assess relevance based on user instructions and past actions. Additionally, Gur et al. (2024) introduce HTML-T5, a specialized language model for web page understanding pre-trained on a large corpus of HTML documents and fine-tuned for extractive summarization. In contrast, we propose training language models to contextualize complex web page observations to enhance decision making in LLM agents.

Automated prompting for closed-source models The quality of outputs from closed-source LLMs heavily depends on the prompts used, leading to a substantial body of research dedicated to prompt engineering to elicit more effective responses from these models. (Kojima et al., 2022;

540 Yao et al., 2022b; Lightman et al., 2024). For example, several recent studies have explored automating the discovery of more effective prompting formats (Shin et al., 2020; Zhou et al., 2023b).
541
542 Mañas et al. (2024) demonstrate that prompt refinement can yield more consistent responses from
543 closed-source models, especially with text-to-image models. Also, Xu et al. (2024) have suggested
544 compressing the input content for LLMs, focusing on the task of retrieval-augmented generation.
545 In this work, we focus on contextualizing the raw observations for agents performing sequential
546 decision making and introduce a novel training algorithm.

547 548 6 CONCLUSION

549
550 In this work, we introduce LCoW, a novel approach to enhancing LLM agents in completing web
551 tasks by leveraging language models to contextualize complex web pages into a more comprehensible
552 form. Our approach separates the understanding of web content from the decision-making
553 process by training a specialized module that generates contextualized representations of complex
554 web pages, which are utilized by LLM agents for enhanced decision making. Through extensive
555 experiments, we demonstrate that this contextualization module significantly improves the decision-
556 making capabilities of LLM agents of varying scales.

557 558 ETHICS STATEMENT

559
560 We introduce LCoW, a framework for improving decision-making capability of LLM agents for
561 web automation. We caution that LLM agents may cause safety issues such as cybersecurity or
562 risks regarding private information, while we believe that LCoW can be valuable for guiding the
563 agents from potential mistakes by providing more contextualized information on the UI elements.
564 Additionally, we believe the improved efficiency and capability of LLM agents with our methods
565 can provide social opportunities to improve user interactions of using digital devices for those with
566 disabilities.

567 568 REPRODUCIBILITY STATEMENT

569
570 For the reproducibility of our results, we have provided a detailed description of our methods and
571 experimental setups in Section 4.1. Additionally, to further facilitate the reproduction, we will
572 release our codes and the model checkpoints in the final version.

573 574 REFERENCES

- 575
576 Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su.
577 Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing*
578 *Systems*, 36, 2024.
- 579
580 Alexandre Drouin, Maxime Gasse, Massimo Caccia, Issam H Laradji, Manuel Del Verme, Tom
581 Marty, Léo Boisvert, Megh Thakkar, Quentin Cappart, David Vazquez, et al. Workarena:
582 How capable are web agents at solving common knowledge work tasks? *arXiv preprint*
583 *arXiv:2403.07718*, 2024.
- 584
585 Hiroki Furuta, Kuang-Huei Lee, Ofir Nachum, Yutaka Matsuo, Aleksandra Faust, Shixiang Shane
586 Gu, and Izzeddin Gur. Multimodal web navigation with instruction-finetuned foundation models.
587 *arXiv preprint arXiv:2305.11854*, 2023.
- 588
589 Izzeddin Gur, Ofir Nachum, Yingjie Miao, Mustafa Safdari, Austin Huang, Aakanksha Chowdhery,
590 Sharan Narang, Noah Fiedel, and Aleksandra Faust. Understanding html with large language
591 models, 2023.
- 592
593 Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and
Aleksandra Faust. A real-world webagent with planning, long context understanding, and program synthesis. In *International Conference on Learning Representations*, 2024.

- 594 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
595 language models are zero-shot reasoners. *Advances in neural information processing systems*,
596 35:22199–22213, 2022.
- 597
- 598 Hanyu Lai, Xiao Liu, Iat Long Iong, Shuntian Yao, Yuxuan Chen, Pengbo Shen, Hao Yu, Hanchen
599 Zhang, Xiaohan Zhang, Yuxiao Dong, et al. Autowebglm: Bootstrap and reinforce a large lan-
600 guage model-based web navigating agent. *arXiv preprint arXiv:2404.03648*, 2024.
- 601 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
602 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *International
603 Conference on Learning Representations*, 2024.
- 604
- 605 Xiao Liu, Tianjie Zhang, Yu Gu, Iat Long Iong, Yifan Xu, Xixuan Song, Shudan Zhang, Hanyu Lai,
606 Xinyi Liu, Hanlin Zhao, et al. Visualagentbench: Towards large multimodal models as visual
607 foundation agents. *arXiv preprint arXiv:2408.06327*, 2024.
- 608 Kaixin Ma, Hongming Zhang, Hongwei Wang, Xiaoman Pan, Wenhao Yu, and Dong Yu. Laser:
609 Llm agent with state-space exploration for web navigation, 2024.
- 610
- 611 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri
612 Alon, Nouha Dziri, Shrimai Prabhunoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad
613 Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine:
614 Iterative refinement with self-feedback, 2023.
- 615 Oscar Mañas, Pietro Astolfi, Melissa Hall, Candace Ross, Jack Urbanek, Adina Williams, Aish-
616 warya Agrawal, Adriana Romero-Soriano, and Michal Drozdal. Improving text-to-image con-
617 sistency via automatic prompt optimization. *arXiv preprint arXiv:2403.17804*, 2024.
- 618
- 619 OpenAI. <https://openai.com/index/hello-gpt-4o/>, 2024.
- 620
- 621 Jiayi Pan, Yichi Zhang, Nicholas Tomlin, Yifei Zhou, Sergey Levine, and Alane Suhr. Autonomous
622 evaluation and refinement of digital agents. In *First Conference on Language Modeling*, 2024.
- 623 Pranav Putta, Edmund Mills, Naman Garg, Sumeet Motwani, Chelsea Finn, Divyansh Garg, and
624 Rafael Rafailov. Agent q: Advanced reasoning and learning for autonomous ai agents. *arXiv
625 preprint arXiv:2408.07199*, 2024.
- 626
- 627 Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Autoprompt:
628 Eliciting knowledge from language models with automatically generated prompts. In *Proceed-
629 ings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
630 Association for Computational Linguistics, 2020.
- 631 Paloma Sodhi, S. R. K. Branavan, Yoav Artzi, and Ryan McDonald. Step: Stacked llm policies for
632 web actions, 2024.
- 633
- 634 Abishek Sridhar, Robert Lo, Frank F. Xu, Hao Zhu, and Shuyan Zhou. Hierarchical prompting
635 assists large language model on web navigation, 2023.
- 636 Zora Zhiruo Wang, Jiayuan Mao, Daniel Fried, and Graham Neubig. Agent workflow memory,
637 2024.
- 638
- 639 Fangyuan Xu, Weijia Shi, and Eunsol Choi. Recomp: Improving retrieval-augmented lms with com-
640 pression and selective augmentation. In *International Conference on Learning Representations*,
641 2024.
- 642
- 643 Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable
644 real-world web interaction with grounded language agents. *Advances in Neural Information Pro-
645 cessing Systems*, 35:20744–20757, 2022a.
- 646 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
647 React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*,
2022b.

648 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,
649 Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for build-
650 ing autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023a.
651
652 Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and
653 Jimmy Ba. Large language models are human-level prompt engineers. In *International Confer-*
654 *ence on Learning Representations*, 2023b.
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A APPENDIX

A.1 WEBARENA-LITE

In this section, we provide experimental results in WebArena-Lite (Liu et al., 2024), consisting of diverse tasks spanning over 6 websites. WebArena-Lite is a benchmark composed of 165 tasks filtered from original WebArena benchmark (Zhou et al., 2023a). Out of entire 812 tasks in WebArena, we consider 165 tasks in WebArena-Lite as an evaluation split, and remaining 647 tasks as a training split. As a seed demonstration, we collect 363 successful trajectories from the training split via utilizing GPT-4o and Claude-3-5-sonnet. As a next step, we train the contextualization module via LCoW, and evaluate it on WebArena-Lite benchmark. As shown in the 3, LCoW results in remarkably better performance compared to the baseline, implying LCoW is effective in learning contextualization in broad ranger of websites.

	GPT-4o
Raw observation	29.7%
LCoW	35.8%

Table 3: We evaluate the success rate of GPT-4o-2024-08-06 and LCoW + GPT-4o-2024-08-06 on WebArena-Lite benchmark. Contextualizing web page observation via LCoW results in improved success rate.

Additionally, we analyze whether LCoW effectively generalizes to task types that were unseen during the training. Specifically, the 812 tasks in the WebArena benchmark are all different but are created based on 190 task templates. Therefore, tasks created from the same task template are different but similar. For example, “What is the top-1 best-selling product in 2022” and “What is the top-3 best-selling product in 2023” are from the same task template. Our research question is whether LCoW can generalize to unseen task templates that are not included in the 363 successful trajectories used for training. Among 165 WebArena-Lite evaluation tasks, 48 tasks belongs to unseen task templates, while 117 tasks belongs to seen task templates. Based on the information, we evaluate success rates on 117 tasks corresponding to seen task templates and 48 tasks corresponding to unseen task templates, respectively. As shown in Table 4, LCoW also generalizes to tasks corresponding to unseen task templates.

	GPT-4o	
	Seen-task-template (117 tasks)	Unseen-task-template (48 tasks)
Raw observation	35.9%	14.5%
LCoW	41.9%	20.8%

Table 4: LCoW also generalizes to unseen tasks, demonstrating more than 6% improvement in success rate on 48 tasks corresponding to the unseen-task-template, as well as 117 tasks corresponding to the seen-task-template.

A.2 GENERALIZATION AT DIFFERENT LEVELS

In this section, we conduct systematic evaluation of LCoW’s generalization capabilities in WorkArena. WorkArena features a two-level task hierarchy: categories at the top level and types within each category. WorkArena provides 7 categories of tasks (i.e., Dashboard, Menus, Service catalog, Knowledge base, Forms, Sort list, and Filter list). Although Sort list and Filter list are sub-category under the List task category, we determine to consider them as a discrete separate task category in order to avoid confusion. For instance, “Form” is a task *category*, and within it, “creating and submitting an incident report” and “creating new user information” are task *types*. We consider two levels of generalization: 1) Unseen-type tasks, tasks of

a different type within the same category (i.e., medium generalization) and 2) Unseen-category tasks, tasks of a different type and category (i.e., hard generalization). Detailed information about seen and unseen tasks are provided in Table 5. In our experiments, we trained the contextualization module on 13 different tasks types and evaluated its performance on 100 individual tasks corresponding to 14 unseen-type tasks and 6 unseen-category tasks. As shown in the Table 5, LCoW demonstrated strong generalization to unseen-type tasks, achieving a 22.6% improvement when using Gemini 1.5-flash as the LLM agent. For example, knowledge learned from tasks of writing change requests and submitting can be generalized to creating new user account. However, we found LCoW struggles to generalize to unseen-category tasks (i.e., *Filter-list*), highlighting the need for greater task diversity in training or enhanced contextual reasoning to address completely new task types.

	GPT-4o		Gemini-1.5-flash	
	Unseen-type	Unseen-category	Unseen-type	Unseen-category
Raw observation	35.7%	0.0%	14.5%	0.0%
LCoW	42.9%	0.0%	37.1%	0.0%

Table 5: LCoW shown to be generalized to unseen task types within the same task category on both GPT-4o and Gemini-1.5-flash backbone, but it struggles to generalize to tasks corresponding to unseen category.

A.3 COMPARISON WITH LLM-BASED PARSER

In this section, we provide experimental results comparing pre-trained LLM-based HTML parser and LCoW. As a LLM-based HTML parser, we utilize Reader-LM-1.5B, an LLM pre-trained to generate concise markdown given complicated HTML. Although Reader-LM generates reasonable summary when given non-lengthy HTML input (i.e., less than 2048 tokens), it tends to generate non-meaningful continuation given HTML longer than 20K tokens, which is prevalent in WorkArena benchmark. Therefore, we utilized accessibility tree observation, a textual web page representation with reduced noisy and enhanced readability. However, Reader-LM tend to repeat the given accessibility tree observations, rather than summarizing or rephrasing them, until it reaches maximum output token limit (i.e., 2048). As a result, as shown in 7 Reader-LM degrades success rates compared to using raw observation directly. Additionally, typical failure cases of Reader-LM are described in 10.

A.4 ENTIRE PROMPTS

In this section, we provide entire prompts used across entire experiments in WorkArena and WebShop.

A.4.1 WORKARENA & WEBARENA

Prompt for contextualization module

```
<system prompt>
You are an agent tasked with extracting and refining a subset of the
webpage's observations based on the content of the page and user
instructions.

<main prompt>
You are currently on the {domain_info} website.
Your task is to generate a "Reasoning" and a "Refined observation"
based on the provided inputs.

First, review the "User instruction" and "History of interactions"
and, then, generate the "Reasoning".
Analyze the progress made so far, and provide a rationale for the
```

Task type	Seen / Unseen
Single-chart-value-retrieval	seen
Single-chart-minmax retrieval	seen
Multi-chart-value retrieval	unseen-type
Multi-chart-minmax retrieval	unseen-type
Create change request	seen
Create problem	seen
Create incident	unseen-type
Create hardware asset	unseen-type
Create user	unseen-type
Knowledge base search	seen
Sort user list	seen
Sort hardware list	seen
Sort asset list	unseen-type
Sort change request list	unseen-type
Sort incident list	unseen-type
Sort service-catalog list	unseen-type
All menu	seen
Impersonation	unseen-type
Order developer laptop	seen
Order iPad mini	seen
Order iPad pro	seen
Order apple watch	seen
Order standard laptop	seen
Order sales laptop	unseen-type
Order apple Macbook pro	unseen-type
Order development laptop PC	unseen-type
Order loander laptop	unseen-type
Filter asset list	unseen-category
Filter change request list	unseen-category
Filter hardware list	unseen-category
Filter incident list	unseen-category
Filter service catalog item list	unseen-category
Filter user list	unseen-category

Table 6: Task type split for systematic evaluation of generalization at different levels in WorkArena benchmark.

	GPT-4o
Raw observation	38.2%
Reader-LM	9.7%
LCoW (iter 1)	44.2%

Table 7: We compare the LCoW and Reader-LM, an LLM pre-trained to summarize web page into markdown format in WorkArena benchmark. However, we observed that Reader-LM struggles to summarize web page observations in WorkArena benchmark.

next steps needed to efficiently accomplish the user instruction on the {domain_info} website.

Second, refine the "AXTree observation at the current time step" into a "Refined observation". Select a subset of the AXTree observation that is essential for completing the user instruction and provide explanations for the corresponding elements in the selected subset.

[Information source]

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

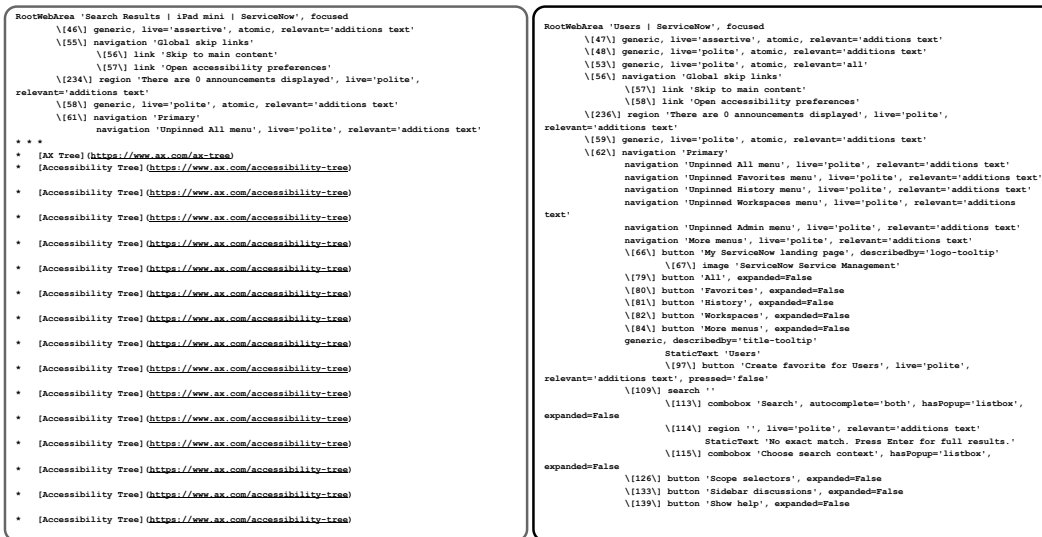


Figure 10: Reader-LM occasionally generates meaningless summarization of given web pages (Left), and mainly re-generates given web pages until the maximum output token limit has reached (Right).

```
# User instruction
{goal}

# History of interactions
{history}

# AXTree observation at the current time step
{observation}
```

Prompt for contextualization module in retry phase

```
<system prompt>
You are an agent tasked with extracting and refining a subset of the
webpage's observations based on the content of the page and user
instructions.

<main prompt>
You are currently on the {domain_info} website.
Your task is to generate a "Reasoning" and a "Refined observation"
based on the provided inputs.

First, review the "User instruction" and "History of interactions"
and, then, generate the "Reasoning".
Analyze the progress made so far, and provide a rationale for the
next steps needed to efficiently accomplish the user instruction on
the {domain_info} website.

Second, refine the "AXTree observation at the current time step"
into a "Refined observation".
Select a subset of the AXTree observation that is necessary for
completing the user instruction.

You may refer to the Hints, which consists of the ground truth next
action, but do not explicitly mention these hints in your output.

[Information source]
```

```

918 # User instruction
919 {goal}
920
921 # History of interactions
922 {history}
923
924 # AXTree observation at the current time step
925 {observation}
926
927 # Hint
928 Ground-truth next action: {action}

```

Prompt for self-contextualization

```

931 <system prompt>
932 You are an agent tasked with extracting and refining a subset of the
933 webpage's observations based on the content of the page and user
934 instructions.
935
936 <main prompt>
937 [General instructions]
938 You are currently on the {domain_info} website.
939 Your task is to generate a "Reasoning" and a "Refined observation"
940 based on the provided inputs.
941
942 First, review the "User instruction" and "History of interactions" and,
943 then, generate the "Reasoning".
944 Analyze the progress made so far, and provide a rationale for the next
945 steps needed to efficiently accomplish the user instruction on
946 the {domain_info} website.
947
948 Second, refine the "AXTree observation at the current time step"
949 into a "Refined observation".
950 Extract a subset of the AXTree observation (e.g., chart, table,
951 menu items) that contains necessary information for completing
952 the user instruction, and explain the extracted elements.
953 Ensure that the information on the elements (e.g., numeric element ID)
954 are correctly included.
955
956 Please follow the format in the [Reasoning & Refinement example]
957 carefully.
958
959 [Information source]
960 # User instruction
961 {goal}
962
963 # History of interactions
964 {history}
965
966 # AXTree observation at the current time step
967 {observation}
968
969 [Reasoning & Refinement example]
970 # Abstract example
971 Here is an abstract version of the answer, describing
the content of each tag.
Make sure you follow this structure, but replace the
content with your own answer:
<reasoning>
Think step by step. Based on the "User instruction,"
"History of interaction," and "AXTree observation at the
current time step":

```

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

```

1. Provide a high-level description of the "AXTree observation at the
current time step."
2. Based on the "User instruction" and "History of interaction,"
track your progress and provide your reasoning on the next action
needed to accomplish the "User instruction."
</reasoning>

<extraction>
Based on your reasoning, identify the elements
(e.g., links, buttons, static text, table row, chart) to focus on.
Then, explain the semantics and functionalities of
each extracted element.
Ensure that:
You do not alter the structure of the AXTree observation.
You extract the element ID (id in [ ]) accurately without any errors.
When extracting chart or table, you must extract the entire chart
or table to avoid any confusion or loss of information.
</extraction>

```

Prompt for LLM agent

```

<system prompt>
You are an agent trying to solve a web task based on the content of
the page and a user instructions.
You can interact with the page and explore.
Each time you submit an action it will be sent to the browser and you
will receive a new page.

<main prompt>
# Instructions
Review the current state of the page and all other information to find
the best possible next action to accomplish your goal.
Your answer will be interpreted and executed by a program, make sure
to follow the formatting instructions.

## Goal:
{goal}

{history}

# Refined observation of current step:
{refined observation}

# Action space:
13 different types of actions are available.
noop(wait_ms: float = 1000)
  Description: Do nothing, and optionally wait for
  the given time (in milliseconds).
  Examples:
    noop()
    noop(500)

send_msg_to_user(text: str)
  Description: Send a message to the user.
  You should send a short answer as a message and
  do not ask questions through message.
  Examples:
    send_msg_to_user('\the city was built in 1751.\')
    send_msg_to_user('\Yes\')
    send_msg_to_user('\No\')
    send_msg_to_user('\31112\')
    send_msg_to_user('\Yoshua Bengio\')

scroll(delta_x: float, delta_y: float)

```

```

1026 Description: Scroll horizontally and vertically.
1027 Amounts in pixels, positive for right or down scrolling,
1028 negative for left or up scrolling. Dispatches a wheel event.
1029 Examples:
1030     scroll(0, 200)
1031     scroll(-50.2, -100.5)
1032
1033 fill(bid: str, value: str)
1034 Description: Fill out a form field.
1035 It focuses the element and triggers an input event
1036 with the entered text. It works for <input>,
1037 <textarea> and [contenteditable] elements.
1038 Examples:
1039     fill('237', 'example value')
1040     fill('45', 'multi-line\nexample')
1041     fill('a12', 'example with "quotes"')
1042
1043 select_option(bid: str, options: str | list[str])
1044 Description: Select one or multiple options in a <select> element.
1045 You can specify option value or label to select.
1046 Multiple options can be selected.
1047 Examples:
1048     select_option('48', 'blue')
1049     select_option('48', ['red', 'green', 'blue'])
1050
1051 click(bid: str, button: Literal['left', 'middle', 'right'] = 'left',
1052 modifiers: list[typing.Literal['Alt', 'Control', 'Meta', 'Shift']]
1053 = [])
1054 Description: Click an element.
1055 Examples:
1056     click('51')
1057     click('b22', button='right')
1058     click('48', button='middle', modifiers=['Shift'])
1059
1060 dblclick(bid: str, button: Literal['left', 'middle', 'right'] =
1061 'left', modifiers: list[typing.Literal['Alt', 'Control', 'Meta',
1062 'Shift']] = [])
1063 Description: Double click an element.
1064 Examples:
1065     dblclick('12')
1066     dblclick('ca42', button='right')
1067     dblclick('178', button='middle', modifiers=['Shift'])
1068
1069 hover(bid: str)
1070 Description: Hover over an element.
1071 Examples:
1072     hover('b8')
1073
1074 press(bid: str, key_comb: str)
1075 Description: Focus the matching element and press a combination of
1076 keys.
1077 It accepts the logical key names that are emitted in the
1078 keyboardEvent.
1079 key property of the keyboard events: Backquote, Minus, Equal,
1080 Backslash, Backspace, Tab, Delete, Escape, ArrowDown, End, Enter,
1081 Home, Insert, PageDown, PageUp, ArrowRight, ArrowUp, F1 - F12,
1082 Digit0 - Digit9, KeyA - KeyZ, etc. You can alternatively specify a
1083 single character you'd like to produce such as "a" or "#".
1084 Following modification shortcuts are also supported: Shift,
1085 Control, Alt, Meta.
1086 Examples:
1087     press('88', 'Backspace')
1088     press('a26', 'Control+a')
1089

```

```

1080         press('a61', 'Meta+Shift+t')
1081
1082     focus(bid: str)
1083         Description: Focus the matching element.
1084         Examples:
1085             focus('b455')
1086
1087     clear(bid: str)
1088         Description: Clear the input field.
1089         Examples:
1090             clear('996')
1091
1092     drag_and_drop(from_bid: str, to_bid: str)
1093         Description: Perform a drag & drop.
1094         Hover the element that will be dragged.
1095         Press left mouse button.
1096         Move mouse to the element that will receive the drop.
1097         Release left mouse button.
1098         Examples:
1099             drag_and_drop('56', '498')
1100
1101     upload_file(bid: str, file: str | list[str])
1102         Description: Click an element and wait for a "filechooser" event,
1103         then select one or multiple input files for upload.
1104         Relative file paths are resolved relative to the current working
1105         directory.
1106         An empty list clears the selected files.
1107         Examples:
1108             upload_file('572', 'my_receipt.pdf')
1109             upload_file('63', ['/home/bob/Documents/image.jpg',
1110                 '/home/bob/Documents/file.zip'])
1111
1112     Only a single action can be provided at once. Example:
1113     fill('a12', 'example with "quotes"')
1114     Multiple actions are meant to be executed sequentially without any
1115     feedback from the page.
1116     Don't execute multiple actions at once if you need feedback from the
1117     page.
1118
1119     # Abstract Example
1120     Here is an abstract version of the answer with description of the
1121     content of each tag.
1122     Make sure you follow this structure, but replace the content with your
1123     answer:
1124
1125     <think>
1126     Think step by step.
1127     If you need to make calculations such as coordinates, write them here.
1128     Describe the effect that your previous action had on the current
1129     content of the page.
1130     </think>
1131
1132     <action>
1133     One single action to be executed.
1134     You can only use one action at a time.
1135     </action>
1136
1137     # Concrete Example
1138     Here is a concrete example of how to format your answer.
1139     Make sure to follow the template with proper tags:
1140
1141     <think>
1142     My memory says that I filled the first name and last name, but I can't

```

```

1134 see any content in the form.
1135 I need to explore different ways to fill the form.
1136 Perhaps the form is not visible yet or some fields are disabled.
1137 I need to replan.
1138 </think>
1139
1140 <action>
1141 fill('a12', 'example with "quotes"')
1142 </action>

```

A.4.2 WEBSHOP

Prompt for contextualization module

```

1147 <system prompt>
1148 You are an agent tasked with extracting and rephrasing a subset of
1149 the webpage's observations based on the content of the page and user
1150 instructions.
1151
1152 <main prompt>
1153 You are currently on the online shopping website.
1154 Your task is to generate a "Reasoning" and a "Refined observation"
1155 based on the provided inputs.
1156
1157 First, review the "User instruction" and "History of interactions"
1158 and, then, generate the "Reasoning".
1159 Analyze the progress made so far, and provide a rationale for the
1160 next steps needed to efficiently accomplish the user instruction on
1161 the online shopping website.
1162
1163 Second, rephrase the "AXTree observation at the current time step"
1164 into a "Rephrased observation".
1165 Select a subset of the AXTree observation that is essential
1166 for completing the user instruction and provide explanations
1167 for the corresponding elements in the selected subset.
1168
1169 [Information source]
1170 # User instruction
1171 {goal}
1172
1173 # History of interactions
1174 {previous_actions}
1175
1176 # AXTree observation at the current time step
1177 {obs}

```

Prompt for self-contextualization

```

1175 <system prompt>
1176 You are an agent tasked with extracting and rephrasing a subset of
1177 the webpage's observations based on the content of the page and user
1178 instructions.
1179
1180 <main prompt>
1181 The current webpage on the web shopping site is described
1182 in the observation.
1183 Evaluate the current progress based on previous actions
1184 and current observation.
1185 Determine the next action by reasoning based on
1186 goal and progress.
1187 Condense the observation into a concise format, highlighting
1188 clickable buttons indicated by [].
1189 Ensure the summary includes only elements relevant to the

```

1188 goal and not already covered in previous actions.
 1189 Focus on clickable buttons indicated as [].
 1190
 1191 Here are a few examples.
 1192
 1193 ****goal****: i would like a 3 ounce bottle of bright citrus
 1194 deodorant for sensitive skin, and price lower than 50.00 dollars
 1195 ****previous actions****:
 1196 1. search[3 ounce bright citrus deodorant sensitive skin]
 1197 ****current observation****:
 1198 [Back to Search]
 1199 Page 1 (Total results: 50)
 1200 [Next >]
 1201 [B078GWRC1J]
 1202 Bright Citrus Deodorant by Earth Mama | Natural and Safe
 1203 for Sensitive Skin, Pregnancy and Breastfeeding,
 1204 Contains Organic Calendula 3-Ounce
 1205 \$10.99
 1206 [B078GTKVXY]
 1207 Ginger Fresh Deodorant by Earth Mama | Natural and Safe for
 1208 Sensitive Skin, Pregnancy and Breastfeeding, Contains Organic
 1209 Calendula 3-Ounce
 1210 \$10.99
 1211 [B08KBVJ4XN]
 1212 Barrel and Oak - Aluminum-Free Deodorant, Deodorant for Men,
 1213 Essential Oil-Based Scent, 24-Hour Odor Protection, Cedar &
 1214 Patchouli Blend, Gentle on Sensitive Skin (Mountain Sage,
 1215 2.7 oz, 2-Pack)
 1216 \$15.95
 1217
 1218 ****rephrased observation****:
 1219 Progress: I searched the keyword '3 ounce bright citrus deodorant
 1220 sensitive skin' to see the relvant items, And now I am looking at
 1221 the item list.
 1222 Reasoning: the next step is to choose an item satisfying the
 1223 specification of bright citrus deodorant.
 1224 I can focus on:
 1225 [B078GWRC1J]
 1226 Bright Citrus Deodorant by Earth Mama | Natural and Safe for
 1227 Sensitive Skin, Pregnancy and Breastfeeding, Contains Organic
 1228 Calendula 3-Ounce
 1229 \$10.99
 1230
 1231 ****goal****: i would like a 3 ounce bottle of bright citrus deodorant
 1232 for sensitive skin, and price lower than 50.00 dollars
 1233 ****previous actions****:
 1234 1. search[3 ounce bright citrus deodorant sensitive skin]
 1235 2. click[B078GWRC1J]
 1236 ****current observation****:
 1237 [Back to Search]
 1238 [< Prev]
 1239 size
 1240 [travel set (4-pack)]
 1241 [3 ounce (pack of 1)]
 [3-ounce (2-pack)]
 scent
 [assorted scents]
 [bright citrus]
 [calming lavender]
 [ginger fresh]
 [simply non-scents]
 Bright Citrus Deodorant by Earth Mama | Natural and Safe for

1242 Sensitive Skin, Pregnancy and Breastfeeding, Contains Organic
 1243 Calendula 3-Ounce
 1244 Price: \$10.99
 1245 Rating: N.A.
 1246 [Description]
 1247 [Features]
 1248 [Reviews]
 1249 [Buy Now]

1250 ****rephrased observation**:**
 1251 Progress: I searched and and clicked the item seems to be most
 1252 relevant to the goal specification. I am looking at the option list.
 1253 Reasoning: As the goal requires 3-ounce bottle, I can focus
 1254 on the size option.
 1255 I can focus on:
 1256 size
 1257 [travel set (4-pack)]
 1258 [3 ounce (pack of 1)]
 1259 [3-ounce (2-pack)]

1260 ****goal**:** i would like a 3 ounce bottle of bright citrus deodorant
 1261 for sensitive skin, and price lower than 50.00 dollars
 1262 ****previous actions**:**
 1263 1. search[3 ounce bright citrus deodorant sensitive skin]
 1264 2. click[B078GWRC1J]
 1265 3. click[3 ounce (pack of 1)]
 1266 ****current observation**:**
 1267 You have clicked 3 ounce (pack of 1).
 1268 [Back to Search]
 1269 [< Prev]
 1270 size
 1271 [travel set (4-pack)]
 1272 [3 ounce (pack of 1)]
 1273 [3-ounce (2-pack)]
 1274 scent
 1275 [assorted scents]
 1276 [bright citrus]
 1277 [calming lavender]
 1278 [ginger fresh]
 1279 [simply non-scents]
 1280 Bright Citrus Deodorant by Earth Mama | Natural and Safe for
 1281 Sensitive Skin, Pregnancy and Breastfeeding, Contains Organic
 1282 Calendula 3-Ounce
 1283 Price: \$10.99
 1284 Rating: N.A.
 1285 [Description]
 1286 [Features]
 1287 [Reviews]
 1288 [Buy Now]

1289 ****rephrased observation**:**
 1290 Progress: I searched and and clicked the item id.
 1291 Among the option list, and I clicked size option.
 1292 Reasoning: According to the progress, I have to focus
 1293 on the scent option as a next step.
 1294 I can focus on:
 1295 scent
 [assorted scents]
 [bright citrus]
 [calming lavender]
 [ginger fresh]
 [simply non-scents]

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

```

**goal**: i would like a 3 ounce bottle of bright citrus
deodorant for sensitive skin, and price lower than 50.00 dollars
**previous actions**:
1. search[3 ounce bright citrus deodorant sensitive skin]
2. click[B078GWRC1J]
3. click[3 ounce (pack of 1)]
4. click[bright citrus]
**current observation**:
You have clicked 3 ounce (pack of 1).
You have clicked bright citrus.
[ Back to Search ]
[ < Prev ]
size
[ travel set (4-pack) ]
[ 3 ounce (pack of 1) ]
[ 3-ounce (2-pack) ]
scent
[ assorted scents ]
[ bright citrus ]
[ calming lavender ]
[ ginger fresh ]
[ simply non-scents ]
Bright Citrus Deodorant by Earth Mama | Natural and Safe
for Sensitive Skin, Pregnancy and Breastfeeding,
Contains Organic Calendula 3-Ounce
Price: $10.99
Rating: N.A.
[ Description ]
[ Features ]
[ Reviews ]
[ Buy Now ]

**rephrased observation**:
Progress: Based on **observation** and **previous actions**,
I clicked size option and scent option.
Reasoning: As there is no more options to select and I met
all requirements specified in the goal, next step is to buy the item.
I can focus on:
[ Buy Now ]

Here is the task.

**goal**:
{goal}
**previous actions**:
{previous_actions}
**current observation**:
{obs}

**rephrased observation**:

```

Prompt for LLM agent

```

Webshop
Instruction:
i would like a 3 ounce bottle of bright citrus deodorant for sensitive
skin, and price lower than 50.00 dollars
[ Search ]

Action: search[3 ounce bright citrus deodorant sensitive skin]

Observation:

```

1350 Progress: I searched the keyword '3 ounce bright citrus deodorant
1351 sensitive skin' to see the relvant items, And now I am looking at the
1352 item list.
1353 Reasoning: Based on the Progress and current observation, the
1354 next step is to choose an item satisfying the specification.
1355 I can focus on:
1356 [B078GWRC1J]
1357 Bright Citrus Deodorant by Earth Mama | Natural and Safe for Sensitive
1358 Skin, Pregnancy and Breastfeeding, Contains Organic Calendula 3-Ounce
1359 \$10.99
1360 [B078GTKVXY]
1361 Ginger Fresh Deodorant by Earth Mama | Natural and Safe for Sensitive
1362 Skin, Pregnancy and Breastfeeding, Contains Organic Calendula 3-Ounce
1363 \$10.99
1364 Action: click[B078GWRC1J]
1365 Observation:
1366 Progress: I searched and and clicked the item seems to be most
1367 relevant to the goal specification.
1368 I am looking at the option list.
1369 Reasoning: As the goal requires 3-ounce bottle, I can focus on
1370 the size option.
1371 I can focus on:
1372 size
1373 [travel set (4-pack)]
1374 [3 ounce (pack of 1)]
1375 [3-ounce (2-pack)]
1376 Action: click[3-ounce (pack of 1)]
1377 Observation:
1378 Progress: I searched and and clicked the item id.
1379 Among the option list, and I clicked size option.
1380 Reasoning: According to the progress, I have to focus on the
1381 scent option as a next step.
1382 I can focus on:
1383 scent
1384 [assorted scents]
1385 [bright citrus]
1386 [calming lavender]
1387 [ginger fresh]
1388 [simply non-scents]
1389 Action: click[bright citrus]
1390 Observation:
1391 Progress: Based on **observation** and **previous actions**, I clicked
1392 size option and scent option.
1393 Reasoning: As there is no more options to select and I met all
1394 requirements specified in the goal, next step is to buy the item.
1395 I can focus on:
1396 [Buy Now]
1397 Action: click[Buy Now]
1398 Observation:
1399 Now Here is the task.
1400 Instruction:
1401 {instruction}
1402
1403

```

1404 {History of observations and actions}
1405
1406 Observation:
1407 {observation}
1408
1409 Action:

```

A.4.3 ACTION MATCHING EVALUATION PROMPT

Prompt for model-based evaluation of action matching

```

1413 <system prompt>
1414 Your task is to evaluate whether the given two action commands are
1415 semantically aligned.
1416
1417 <main prompt>
1418 You will be given
1419 1). reference action which indicates an correct action.
1420 2). predicted action which is predicted by assistant agent
1421
1422 Your task is to assess whether the message in predicted action
1423 is semantically aligned with message in the reference action.
1424 Please make sure you read and understand these instructions
1425 carefully.
1426 Please keep this document open while reviewing, and refer to it as
1427 needed.
1428
1429 Evaluation Criteria:
1430 Alignment = 1: the predicted action is semantically aligned with the
1431 reference action.
1432     send_msg_to_user('30%') and send_msg_to_user('The percentage of
1433     amount of pending orders among entire orders is 30%') are
1434     semantically aligned.
1435     click('a34') and click('a34', button='left') is semantically
1436     aligned.
1437
1438 Alignment = 0: the predicted action is semantically not aligned with
1439 the reference action.
1440     send_msg_to_user('$25') and send_msg_to_user('The requested
1441     value is $29') are not semantically aligned.
1442     click('a34') and click('a34', button='left') are semantically
1443     aligned.
1444
1445 Evaluation Steps:
1446 1. Write a simple feedback that assess whether the predicted action
1447 is semantically aligned with the reference action.
1448 2. After writing a feedback, write a score that is 0 or 1. You
1449 should refer to the Evaluation Criteria.
1450 3. The output format should look as follows: "Feedback: (write a
1451 feedback for criteria) [RESULT] (an integer number among 0 or 1)"
1452 4. Please do not generate any other opening, closing, and
1453 explanations.
1454
1455 reference action: {ref_action}
1456 predicted action: {pred_action}
1457
1458 Feedback:

```

A.5 QUALITATIVE ANALYSIS

In this section, we provide 1) concrete examples of webpages contextualized by LCoW and 2) examples of action sequence of LLM agents with and without LCoW.

1458 A.5.1 WEB PAGE CONTEXTUALIZATION
1459

1460 As demonstrated in the following examples, contextualization model trained via LCoW provides
1461 contexts (e.g., role of the UI elements specific to the website) that enables more accurate decision
1462 making of LLM-based web agent.

1463 **Example 1**
1464

1465 **Instruction:** Sort the “users” list by the following fields:- Updates (descending)
1466

1467 **Contextualized web page observation:**

1468 [a183] grid 'Users.', multiselectable=False
1469 This is the main table containing the list of users.
1470 We need to focus on this element to find ways to interact
1471 with the column header.
1472 [a186] rowgroup ''
1473 [a187] row ''
1474 [a188] columnheader 'Select All'
1475 [a193] columnheader '\uf1e4 Show column search row'
1476 [a197] columnheader 'User ID User ID column options'
1477 [a204] columnheader 'Name Name column options'
1478 [a211] columnheader 'Email Email column options'
1479 [a218] columnheader 'Avatar Avatar column options'
1480 [a225] columnheader 'Title Title column options'
1481 [a232] columnheader 'Company Company column options'
1482 [a239] columnheader 'Department Department column options'
1483 [a246] columnheader 'Location Location column options'
1484 [a253] columnheader 'Time zone Time zone column options'
1485 [a260] columnheader 'Business phone Business phone column options'
1486 These are the column headers in the users table.
1487 We need to focus on one of these elements to trigger
1488 the sorting action.
1489 However, the AXTree does not show any direct options to
1490 sort these columns.
1491
1492 [a77] button '\uf13e Personalize List'
1493 This button might provide options to customize
1494 the list view, potentially including sorting options.

1491 **Example 2**
1492

1493 **Instruction:** Go to the hardware store and order 10 “Developer Laptop (Mac)” with configuration ‘Ad-
1494 ditional software requirements’: ‘Trello, Zoom, Microsoft Office 365, HubSpot, Adobe Creative Cloud’,
1495 ‘Adobe Acrobat’: True, ‘Adobe Photoshop’: True, ‘Eclipse IDE’: True
1496

1497 **Contextualized web page observation:**

1498 [a184] group ''
1499 [a187] Legend ''
1500 [a189] heading 'Optional Software'
1501 [a194] LayoutTable ''
1502 [a195] checkbox 'Adobe Acrobat', checked='true'
1503 [a196] LabelText ''
1504 [a201] LayoutTable ''
1505 [a202] checkbox 'Adobe Photoshop', checked='false'
1506 [a203] LabelText ''
1507 [a208] LayoutTable ''
1508 [a209] checkbox 'Eclipse IDE', checked='true'
1509 [a210] LabelText ''
1510 This section contains checkboxes for optional software.
1511 The "Adobe Acrobat" checkbox [a195] is already checked,
which aligns with the user’s requirements.

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521

The "Adobe Photoshop" checkbox [a202] is currently unchecked, and we need to check it to meet the user's specifications.

The "Eclipse IDE" checkbox [a209] is already checked, which is correct according to the user's requirements.

The LabelText elements [a196], [a203], and [a210] associated with each checkbox are likely the clickable areas to toggle the checkbox state. We should interact with the LabelText for "Adobe Photoshop" to check the box.

1522
1523
1524

Example 3

Instruction: What is the value of "broken" in the chart (in percent)?

Contextualized web page observation:

we should focus on the pie chart elements.
The relevant section is:

1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552

```
[a962] region 'Pie on Change Request table by Reason
created by Henry Turner: 2024-09-26 05:41:16.
Highcharts interactive chart.'
[a963] region 'Chart screen reader information,
Pie on Change Request table by Reason created by
Henry Turner:
2024-09-26 05:41:16.'
StaticText 'Pie chart with 2 slices.'
generic '', live='polite', relevant='additions text'
generic '', live='assertive', relevant='additions text'
generic '', live='polite', relevant='additions text'
[a973] SvgRoot 'Interactive chart'
group ''
group ''
image '1. (empty), 30.'
image '2. broken, 3.'
[a977] region 'Toggle series visibility,
Pie on Change Request table by Reason created by Henry
Turner: 2024-09-26 05:41:16'
[a978] list ''
[a979] listitem ''
[a980] button 'Show (empty)', pressed='true'
[a981] listitem ''
[a982] button 'Show broken', pressed='true'
[a984] button 'View chart menu, Pie on Change Request
table by Reason created by Henry Turner:
2024-09-26 05:41:16', expanded=False
StaticText 'End of interactive chart.'
```

This section contains the pie chart. The chart has two slices: "(empty)" with a value of 30, and "broken" with a value of 3. The "broken" slice is the one we need to focus on to answer the user's question.

1553
1554
1555
1556
1557

Example 4

Instruction: Sort the "users" list by the following fields: - Language (ascending)

Contextualized web page observation:

we should focus on the following elements:

1561
1562
1563
1564
1565

- [a1053] button '\uf135 Add selected options to the Selected listbox' This button is used to add the selected Language field to

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

the Selected list in the personalization dialog.

2. [a1108] button 'OK'
This button is used to close the personalization dialog and apply the changes.
3. [a246] columnheader 'Location Location column options'
[a248] button 'Location'
[a252] button 'Location column options', hasPopup='menu'
These elements represent the sorting control for the Location column, which might include the Language information.
We'll need to interact with this to sort the list after closing the personalization dialog.

Note: The Language column is not directly visible in the current AXTree observation, likely because it hasn't been added to the list yet.

After adding it and closing the dialog, we would need to find the newly added Language column header to sort it.

Example 5

Instruction: Create a new user with a value of "8273378363827782819" for field "User ID", a value of "Webster-Richard" for field "Last name", a value of "Lynn-David" for field "First name", a value of "lynn-david.webster-richard@workarena.com" for field "Email", a value of "false" for field "Internal Integration User", a value of "Outlook" for field "Calendar integration", a value of "(514) 357-8997" for field "Mobile phone", a value of "false" for field "Password needs reset", a value of "English" for field "Language", and a value of "false" for field "Web service access only".

Contextualized web page observation:

we should focus on the following elements:

```
[a298] LabelText ''
      [a299] note ''
      StaticText 'Internal Integration User'
[a305] checkbox 'Internal Integration User', checked='false'
```

This element represents the "Internal Integration User" field. It's a checkbox that is currently unchecked (false), which aligns with the user instruction. No action is needed for this field.

```
[a353] LabelText ''
      [a354] note ''
      StaticText 'Calendar integration'
[a359] combobox 'Calendar integration' value='Outlook',
hasPopup='menu', expanded=False
      [a360] menuitem '-- None --', selected=False
      [a361] menuitem 'Outlook', selected=True
```

This element represents the "Calendar integration" field. It's a combobox that is already set to "Outlook", which aligns with the user instruction. No action is needed for this field.

```
[a429] LabelText ''
      [a430] note ''
      StaticText 'Mobile phone'
[a435] textbox 'Mobile phone'
```

This element represents the "Mobile phone" field.

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

```

It's a textbox where we need to enter "(514) 357-8997" as per
the user instruction.

[a327] LabelText ''
  [a328] note ''
  StaticText 'Language'
[a333] combobox 'Language' value='-- None --',
      hasPopup='menu', expanded=False
  [a334] menuitem '-- None --', selected=True
  [a335] menuitem 'English', selected=False

This element represents the "Language" field.
It's a combobox where we need to select "English" as per the
user instruction.

[a285] LabelText ''
  [a286] note ''
  StaticText 'Web service access only'
[a292] checkbox 'Web service access only', checked='false'

```

A.5.2 ACTION SEQUENCE

In this subsection, we provide the action sequences obtained by rolling out the LLM agent on evaluation tasks in Figure 11. For a given instruction, we compare the rollouts with and without the application of LCoW, highlighting the differences in decision-making behavior and task completion efficiency.

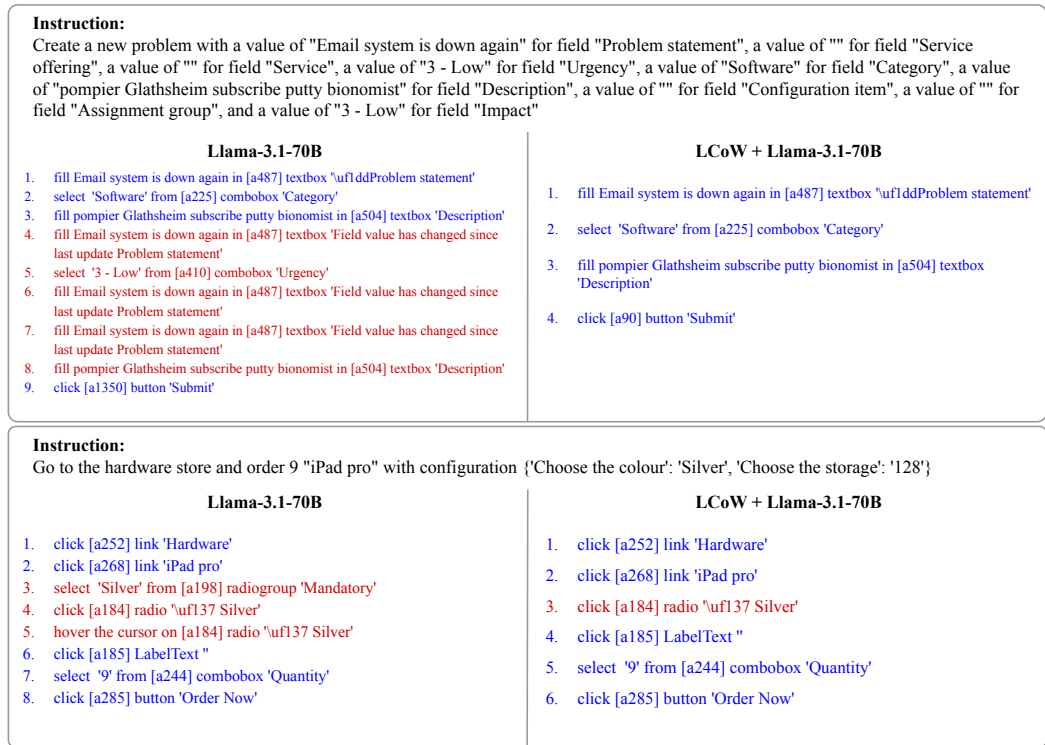


Figure 11: Llama-3.1-70B agent making decisions based on raw observations (**left**) completes the task but exhibits inefficient behavior, such as repeating actions it has already performed or deciding on redundant actions, which are indicated as a red color. In contrast, the Llama-3.1-70B agent making decisions based on contextualized observations via LCoW (**right**) completes the task efficiently, with few or no mistakes and minimal repetition of redundant actions.

1674 A.6 EXPERIMENTAL DETAILS
 1675

1676	Task type	Number of seed demonstrations
1677		
1678	Multi-chart-value retrieval	8
1679	Multi-chart-minmax retrieval	12
1680	Single-chart-value-retrieval	12
1681	Single-chart-minmax retrieval	11
1682	Create change request	8
1683	Create incident	0
1684	Create hardware asset	0
1685	Create problem	14
1686	Create user	11
1687	Knowledge base search	12
1688	Filter asset list	0
1689	Filter change request list	0
1690	Filter hardware list	0
1691	Filter incident list	0
1692	Filter service catalog item list	0
1693	Filter user list	0
1694	Sort asset list	0
1695	Sort change request list	0
1696	Sort hardware list	3
1697	Sort incident list	3
1698	Sort service-catalog list	5
1699	Sort user list	4
1700	All menu	15
1701	Impersonation	12
1702	Order developer laptop	15
1703	Order iPad mini	14
1704	Order iPad pro	15
1705	Order sales laptop	15
1706	Order standard laptop	15
1707	Order apple watch	15
1708	Order apple Macbook pro	15
1709	Order development laptop PC	15
1710	Order loader laptop	15
1711	Total	264

1707
 1708 **Table 8: We collected trajectories from 15 individual tasks for each 33 task types in WorkArena**
 1709 **benchmark using GPT-4o-0806 and Claude-3.5-sonnet agent, thereby collecting 264 successful tra-**
 1710 **jectories. We utilized them as a seed demonstrations for LCoW.**

1711
 1712
 1713
 1714
 1715
 1716
 1717
 1718
 1719
 1720
 1721
 1722
 1723
 1724
 1725
 1726
 1727