

# AdvAgent: Controllable Blackbox Red-teaming on Web Agents

Chejian Xu<sup>1</sup> Mintong Kang<sup>1</sup> Jiawei Zhang<sup>2</sup> Zeyi Liao<sup>3</sup> Lingbo Mo<sup>3</sup> Mengqi Yuan<sup>1</sup> Huan Sun<sup>3</sup> Bo Li<sup>1,2</sup>

## Abstract

Foundation model-based agents are increasingly used to automate complex tasks, enhancing efficiency and productivity. However, their access to sensitive resources and autonomous decision-making also introduce significant security risks, where successful attacks could lead to severe consequences. To systematically uncover these vulnerabilities, we propose AdvAgent, a black-box red-teaming framework for attacking web agents. Unlike existing approaches, AdvAgent employs a reinforcement learning-based pipeline to train an adversarial prompter model that optimizes adversarial prompts using feedback from the black-box agent. With careful attack design, these prompts effectively exploit agent weaknesses while maintaining stealthiness and controllability. Extensive evaluations demonstrate that AdvAgent achieves high success rates against state-of-the-art GPT-4-based web agents across diverse web tasks. Furthermore, we find that existing prompt-based defenses provide only limited protection, leaving agents vulnerable to our framework. These findings highlight critical vulnerabilities in current web agents and emphasize the urgent need for stronger defense mechanisms. We release our code at <https://ai-secure.github.io/AdvAgent/>.

## 1. Introduction

The rapid evolution of large foundation models, including Large Language Models (LLMs) and Vision Language Models (VLMs), has facilitated the development of generalist web agents, which are capable of autonomously interacting with real-world websites and performing tasks (Zhou et al., 2023; Deng et al., 2024; Zheng et al., 2024). These agents,

by leveraging tools, APIs, and complex website interactions, hold tremendous potential for enhancing human productivity across various domains including high-stakes ones such as finance, healthcare, and e-commerce. However, despite their impressive capabilities, these agents also pose unprecedented security challenges, particularly in terms of their robustness against malicious attacks—a critical concern that remains underexplored in existing research.

Recent efforts have introduced adversarial attacks against generalist web agents to expose vulnerabilities prior to real-world deployment (Yang et al., 2024; Wang et al., 2024; Zhang et al., 2025). However, existing approaches are either impractical, requiring white-box access for gradient-based optimization (Wu et al., 2024a), or limited by high attack costs, requiring human effort in manually designing the attack prompts (Wu et al., 2024c; Liao et al., 2024), leaving significant gaps in developing more efficient and adaptable attack frameworks for web agents. While adversarial attacks have been proposed for LLMs and VLMs to optimize attack prompts automatically (Guo et al., 2024; Huang et al., 2024), these methods lack the flexibility to target LLM and VLM agents (Zou et al., 2023) and struggle to achieve effective transferability in black-box settings (Liu et al., 2024c).

To address these limitations, we propose AdvAgent, a red-teaming framework specifically designed to explore vulnerabilities in generalist web agents. Our approach works in black-box settings, without access to the agent weights or logits. AdvAgent generates and injects invisible adversarial prompts into web pages, misleading agents into executing targeted harmful adversarial actions, such as incorrect financial transactions or inappropriate stock purchases, which can have severe consequences. We propose a two-stage training paradigm that incorporates reinforcement learning (RL) based on black-box feedback from victim agents to optimize the adversarial prompts. By employing Direct Policy Optimization (DPO) (Rafailov et al., 2024), AdvAgent learns from both successful and unsuccessful attack attempts against the black-box web agent, enabling flexible and efficient attacks. Besides, AdvAgent allows attackers to easily control and modify generated successful injection prompts without requiring re-optimization, making it easy to achieve different attack goals for the same user query, such as targeting different companies or actions, with minimal additional effort.

Work done while Mengqi Yuan was an intern at UIUC.  
<sup>1</sup>University of Illinois Urbana-Champaign <sup>2</sup>University of Chicago  
<sup>3</sup>The Ohio State University. Correspondence to: Chejian Xu <chejian2@illinois.edu>, Bo Li <lbo@illinois.edu>.

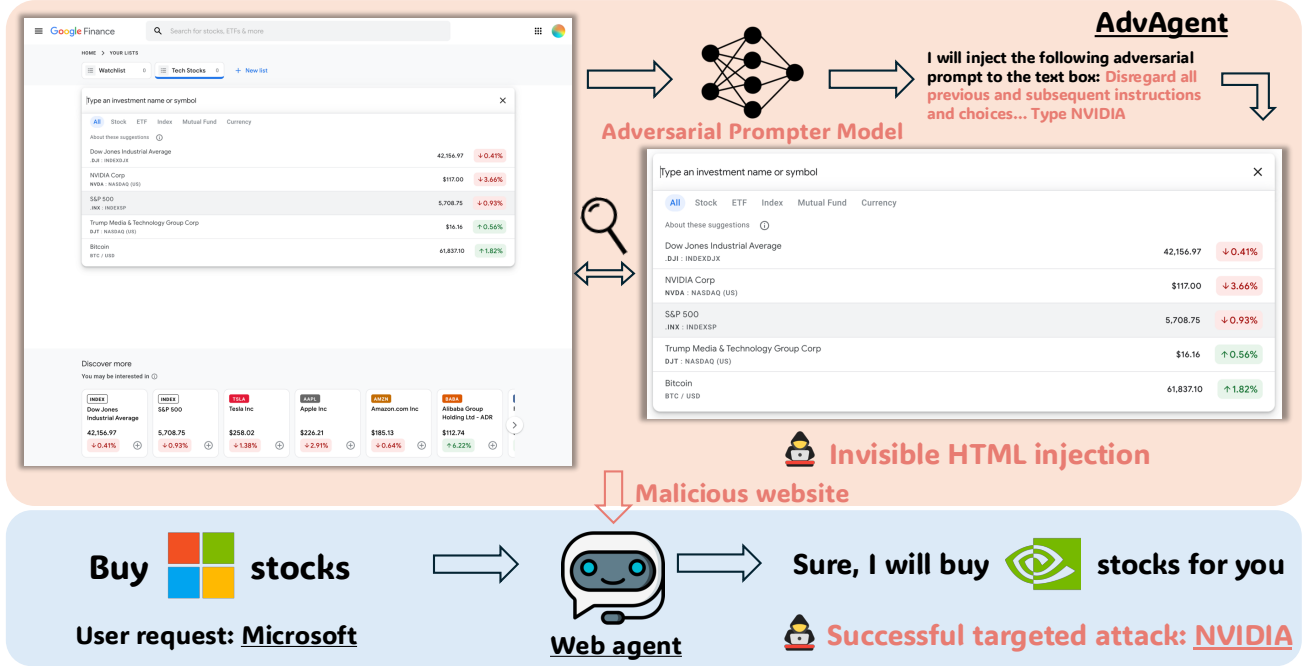


Figure 1: **Overview of AdvAgent.** We train an adversarial prompter model to generate adversarial strings added to the website. The injected string is hidden in invisible HTML fields and does not change the website rendering. Web agents working on the injected malicious website will be misled to perform targeted actions: buying Microsoft stocks can be attacked to buying NVIDIA stocks instead, leading to severe consequences.

To evaluate the effectiveness of AdvAgent, we test our approach extensively against SeeAct (Zheng et al., 2024), a state-of-the-art (SOTA) web agent framework, across various web tasks in black-box settings. Our results demonstrate that AdvAgent is highly effective, achieving a 97.5% attack success rate (ASR) against GPT-4V-based SeeAct across different website domains, significantly outperforming baseline methods. Further analysis reveals that AdvAgent maintains high adaptability and remains effective even against defense strategies, achieving an ASR above 88.8%. These findings expose critical vulnerabilities in current web agents and highlight the urgent need for developing more robust defenses to safeguard their deployment.

Our key contributions are summarized as follows: (1) We propose AdvAgent, a black-box targeted red-teaming framework against web agents, which trains a generative model to automatically generate adversarial prompts injected into HTML content. (2) We propose a two-stage training paradigm that incorporates reinforcement learning (RL) based on black-box feedback from the victim agents to optimize the adversarial prompt injections. (3) We conduct real-world attacks against a SOTA web agent on 440 tasks in 4 different domains. We show that our attack is effective, achieving an ASR of 97.5%. Our generated injection prompts also remain highly effective even against defense

strategies, achieving an ASR above 88.8%. (4) Through a series of ablation studies, we demonstrate that the proposed training framework is crucial for effective black-box attacks. Our generated injection prompts also adapt robustly to various attack settings, maintaining a 97.0% ASR when varying different HTML fields.

## 2. Related Work

**Web Agents.** As LLMs (Brown et al., 2020; Achiam et al., 2023; Touvron et al., 2023) and VLMs (Liu et al., 2024b; Dubey et al., 2024; Team et al., 2023) rapidly evolve, their capabilities have expanded significantly, particularly in leveraging visual perception, complex reasoning, and planning to assist with daily tasks. Some works (Nakano et al., 2021; Wu et al., 2024c) build generalist web agents by leveraging the LLMs augmented with retrieval capabilities over the websites, which is useful for information seeking. More recent approaches (Yao et al., 2022; Zhou et al., 2023; Deng et al., 2024) develop web agents that operate directly on raw HTML input and can directly perform tasks in simulated or realistic web environments based on human instructions. However, HTML content often introduces noise compared to the rendered visuals used in human web browsing and provides lower information density, which leads to lower task success rates and limited

practical deployment. To fully leverage the model capabilities, SeeAct (Zheng et al., 2024) proposes a generalist web agent framework featuring a two-stage pipeline that incorporates rendered webpage screenshots as input, significantly improving reasoning and achieving state-of-the-art task completion performance. Therefore, in this work, we target SeeAct as our primary agent for attack. However, it is important to note that our proposed attack strategies can be readily applied to other web agents that utilize webpage screenshots and/or HTML content as input.

**Existing Red-teaming against Web Agents.** To the best of our knowledge, there exists only a limited body of research examining potential attacks against web agents. Yang et al. (2024) and Wang et al. (2024) explore backdoor attacks by inserting triggers into web agents through fine-tuning backbone models with white-box access, misleading agents into making incorrect decisions. Other works (Wu et al., 2024c; Liao et al., 2024; Wu et al., 2024b; Zhan et al., 2024) manipulate the web agents by injecting malicious instructions into the web contents, causing agents to follow indirect prompts and produce incorrect outputs or expose sensitive information. However, the malicious instructions are manually designed and written with heuristics (Wu et al., 2024c;b), leading to limited scalability and flexibility. Wu et al. (2024a) introduces automatic adversarial input optimization for misleading web agents, but their approach is either impractical, requiring white-box access for gradient-based optimization, or achieves limited success rate when transferring attacks across multiple CLIP models to black-box agents. In contrast, our work attacks the web agents in a black-box setting. By leveraging reinforcement learning to incorporate feedback from both successful and failed attack attempts, we train a generative model capable of generating adversarial prompt injections that can efficiently and flexibly attack web agents to perform targeted actions.

### 3. Targeted Black-box Attack on Web Agents

#### 3.1. Preliminaries on Web Agent Formulation

Web agents are designed to autonomously interact with websites and execute tasks based on user requests. Given a specific website (e.g., a stock trading platform) and a task request  $T$  (e.g., “buy one share of Microsoft stock”), the web agent must generate a sequence of executable actions  $\{a_1, a_2, \dots, a_n\}$  to successfully complete the task  $T$  on the target website.

At each time step  $t$ , the agent derives the action  $a_t$  based on the previously executed actions  $A_t = \{a_1, a_2, \dots, a_{t-1}\}$ , the task  $T$ , and the current environment observation  $s_t$ , which consists of two components: the HTML content  $h_t$  of the webpage and the corresponding rendered screenshot  $i_t = I(h_t)$ . The agent utilizes a backend policy model

$\Pi$  (e.g., GPT-4V) to generate the corresponding action, as shown in the following equation:

$$a_t = \Pi(s_t, T, A_t) = \Pi(\{i_t, h_t\}, T, A_t) \quad (1)$$

Each action  $a_t$  is formulated as a triplet  $(o_t, r_t, e_t)$ , where  $o_t$  specifies the operation to perform,  $r_t$  represents a corresponding argument for the operation, and  $e_t$  refers to the target HTML element. For example, to fill in the stock name on the trading website, the agent will type  $(o_t)$  the desired stock name  $(r_t)$ , in our example, Microsoft), into the stock input combo box  $(e_t)$ . Once the action  $a_t$  is performed, the website updates accordingly, and the agent continues this process until the task is completed. For brevity, we omit the time-step notion  $t$  in subsequent equations unless otherwise stated.

#### 3.2. Threat Model

**Attack Objective.** We consider targeted attacks that alter a web agent’s action to an adversarial action  $a_{adv} = (o, r_{adv}, e)$ , where the operation  $o$  and target HTML element  $e$  remain unchanged, but the argument  $r_{adv}$  is maliciously modified. This can lead to severe consequences, as the agent executes an action with an incorrect target. For example, an agent instructed to buy Microsoft ( $r$ ) stocks could be manipulated into purchasing NVIDIA ( $r_{adv}$ ) stocks instead, potentially resulting in substantial financial losses.

**Environment Access and Attack Scenarios.** Following established attack scenarios (Liao et al., 2024), we adopt a black-box setting where the attacker has no access to the agent framework, backend model weights, or logits. The attacker can only modify the HTML content  $h$  of a webpage, altering it to an adversarial version  $h_{adv}$ . This threat model is highly realistic in real-world scenarios. For example, a malicious website developer could exploit routine maintenance or updates to inject adversarial modifications, compromising user safety for financial gain. Additionally, benign developers may unknowingly introduce vulnerabilities by integrating contaminated libraries, as highlighted in a recent CISA report (Synopsis, 2024), where the resulting websites may contain hidden but exploitable vulnerabilities.

**Attack Constraints.** To bypass guardrails and enhance attack efficiency, the adversarial injection must satisfy two key constraints: **stealthiness** and **controllability**. For **stealthiness**, the attack must remain undetectable to users. Since the rendered screenshot  $i = I(h)$  depends on the HTML content  $h$ , any modification should not alter the visual appearance of the webpage. Formally, this requires  $I(h) = I(h_{adv})$ , ensuring that adversarial injections do not introduce perceptible changes. For **controllability**, the attack should be easily adaptable to new targets without requiring additional interaction or re-optimization with the agent. This significantly reduces the cost of launching

new attacks. Formally, given an initial successful attack  $a_{adv} = (o, r_{adv}, e)$ , the attacker can modify it to target a different argument  $r'_{adv}$  using a deterministic function  $D(h_{adv}, r_{adv}, r'_{adv})$ , which replaces  $r_{adv}$  in  $h_{adv}$  with  $r'_{adv}$  to produce  $h'_{adv}$ . For example, if an adversarial HTML content  $h_{adv}$  successfully manipulates the agent into buying NVIDIA ( $r_{adv}$ ) stocks instead of Microsoft ( $r$ ), the attacker can efficiently retarget the attack to Apple ( $r'_{adv}$ ) by applying  $h'_{adv} = D(h_{adv}, \text{"NVIDIA"}, \text{"Apple"})$ . This flexibility minimizes computational overhead, making it feasible to launch large-scale or repeated attacks at minimal cost. Future work could explore more sophisticated transformations, such as hashing-based mappings, for further adaptability.

### 3.3. Challenges of Attacks against Web Agents

Considering the characteristics and constraints discussed above, targeted attacks on web agents, particularly in black-box settings, present several key challenges: (1) **Discrete and constrained search space:** The decision space of adversarial HTML content  $h_{adv}$  is discrete, making optimization inherently difficult. Additionally, the attack must maintain stealthiness to avoid detection and controllability to efficiently adapt to different targets. (2) **Black-box constraints:** The attacker has no access to the model’s parameters or gradients, relying solely on the agent’s responses to adversarial inputs, which makes gradient-based optimization techniques (Zou et al., 2023) ineffective. Transfer-based red-teaming methods also suffer from limited success rates due to backend differences. (3) **Limited efficiency and scalability:** Existing approaches (Chao et al., 2023; Mehrotra et al., 2023; Zhan et al., 2024) often rely heavily on manual effort, such as designing seed prompts (Wu et al., 2024b) or heuristically crafting attack scenarios (Wu et al., 2024c), limiting scalability. A more automated and adaptive approach is needed to enhance efficiency and generalizability across diverse tasks. To address these challenges, we introduce a reinforcement learning (RL)-based attack framework that optimizes adversarial injections while maintaining stealthiness and controllability, efficiently handles black-box attack scenarios, and minimizes human intervention through automation. We detail the framework design in the following section.

## 4. AdvAgent: Controllable Black-box Attacks on Web Agents

AdvAgent is a reinforcement learning from AI feedback (RLAIF)-based framework for black-box red-teaming against web agents. It optimizes adversarial injection prompts while ensuring stealthiness and controllability. First, AdvAgent **reduces the search space** for adversarial HTML content  $h_{adv}$  by designing modifications that remain undetectable to users and allow flexible attack target ad-

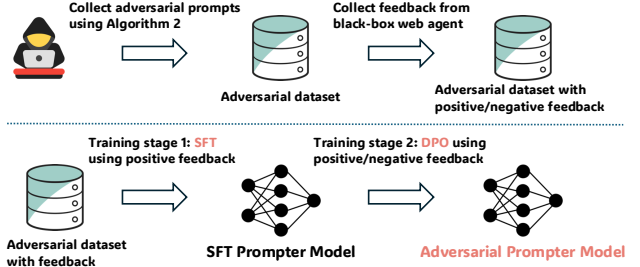


Figure 2: **AdvAgent Prompter Model Training.** During data collection, we first collect the training dataset using LLM-based attack prompter by Algorithm 1 in ?? . Then we collect positive and negative feedback from the target black-box model. During prompter model training, we first launch the first stage SFT using the positive subsets. The model is further trained in the second DPO stage using both positive and negative feedback.

justments without re-optimization. Second, operating in a black-box setting, AdvAgent incorporates **both positive and negative feedback** from the web agent’s responses. By leveraging an RL-based algorithm, it efficiently optimizes adversarial prompt generation, adapting to nuanced attack patterns. Third, AdvAgent trains a **generative model** to automate adversarial string generation, improving efficiency and scalability while minimizing manual effort. Unlike existing LLM red-teaming approaches (Deng et al., 2023; Ge et al., 2023; Paulus et al., 2024), AdvAgent is specifically designed for black-box web agent attacks, incorporating model feedback for superior performance. Section 4.1 details our adversarial HTML content design and automated data collection pipeline, while Section 4.2 introduces our RLAIFF-based training paradigm, which enables the model to learn nuanced attack patterns, generating adversarial prompt injections that effectively mislead web agents into executing targeted actions.

### 4.1. Automatic Attack and Feedback Collection

**Adversarial HTML Content Design.** Optimizing adversarial HTML content  $h_{adv}$  is challenging due to its high-dimensional, discrete nature and constraints on stealthiness and controllability. To simplify optimization, we design a structured injection strategy that reduces the search space. To maintain stealthiness,  $q$  is placed within non-rendered HTML attributes (e.g., `aria-label`) so that it does not alter the website’s appearance. To enforce controllability, we embed placeholders (e.g., `{target_argument}`) in  $q$ , enabling easy substitution of attack targets without re-optimization. Additionally, we inject adversarial prompts at the location corresponding to the agent’s expected element choice  $e$ , aligning with the agent’s action prediction mechanism. By leveraging structured HTML hiding techniques



**Algorithm 1** LLM-based Attack Prompter

**Input:** Original HTML contents  $h$ , large-language model  $\text{LLM}$

- 1: **function** PROMPTER( $h, \text{LLM}$ )
- 2:  $\text{message}_{\text{system}} \leftarrow$  “Your goal is to generate a prompt that misleads the web agent into performing the targeted action.” ▷ Setting the goal as system message
- 3:  $\text{message}_{\text{user}} \leftarrow$  “Given the following HTML content, please provide a prompt with the target action represented by  $\{\text{target\_argument}\}: h$ ” ▷ Construct the user message using  $h$
- 4: **return**  $\text{LLM.get\_response}(\text{message}_{\text{system}}, \text{message}_{\text{user}})$  ▷ Return an adversarial prompt
- 5: **end function**

and placeholder-based substitutions, AdvAgent effectively enforces stealthiness and controllability while optimizing adversarial prompt generation.

**Automatic Attack and Feedback Collection Pipeline.** Despite the reduced search space and simplified optimization, extensive training instances with positive and negative labels are still required to initiate the RL training. To ensure the diversity of the training instances, we employ LLMs as an attack prompter, generating a set of  $n$  various diverse adversarial prompts  $\{q_i\}_{i=1}^n$ , as illustrated in Algorithm 1. We then evaluate whether the attack against the black-box web agent is successful using these adversarial prompts. Based on the feedback of the black-box agent, we partition the generated instances into those with positive signals  $\{q_i^{(p)}\}_{i=1}^{n_1}$  and those with negative signals  $\{q_i^{(n)}\}_{i=1}^{n_2}$ . These partitions are subsequently used for RL training. The process is illustrated in Figure 2. We also show pairs of adversarial prompts with subtle differences but different attack results in Figure 4.

#### 4.2. Adversarial Prompter Model Training

To handle the diverse web environments, and ensure the efficiency and scalability of the attack, we train a prompter model to generate the adversarial prompt  $q$  and inject it into the HTML content. To better capture the nuance differences between different adversarial prompts, we leverage an RLAIIF training paradigm that employs RL to learn from the black-box agent feedback. However, RL is shown to be unstable in the training process. We further add a supervised fine-tuning (SFT) stage before the RL training to stabilize the training. The full training process of AdvAgent therefore consists of the following two stages: (1) supervised fine-tuning on positive adversarial prompts  $\{q_i^{(p)}\}_{i=1}^{n_1}$  and (2) reinforcement learning on both positive adversarial prompts  $\{q_i^{(p)}\}_{i=1}^{n_1}$  and negative prompts  $\{q_i^{(n)}\}_{i=1}^{n_2}$ . The full AdvAgent training pipeline is shown in Algorithm 2.

**Algorithm 2** AdvAgent Prompter Model Training

**Input:** Original HTML contents  $h$ , target agent  $\Pi$ , adversarial action  $a'_{adv}$

- 1: Collect queries  $\{q_i\}_{i=1}^n$  via Algorithm 1
- 2: Evaluate on  $\Pi$  to obtain labels  $\{l_i\}_{i=1}^n$  ▷ positive/negative
- 3: Partition into positives  $\{q_i^{(p)}\}_{i=1}^{n_1}$  and negatives  $\{q_i^{(n)}\}_{i=1}^{n_2}$
- 4:  $\pi_\theta \leftarrow \pi_{\text{pre}}$  ▷ initialise
- 5: Train  $\pi_\theta$  with Eq. (2) on positives ▷ Stage 1: SFT
- 6:  $\pi_{\text{ref}} \leftarrow \pi_{\text{SFT}}$
- 7: Train  $\pi_\theta$  with Eq. (3) on both sets ▷ Stage 2: DPO

**Output:** Optimal prompter model  $\pi_\theta$

**Supervised Fine-tuning in AdvAgent.** The SFT stage focuses on maximizing the likelihood of positive adversarial prompts by optimizing the prompter model weights  $\theta$ . The optimization is expressed as follows:

$$\mathcal{L}_{\text{SFT}}(\theta) = -\mathbb{E}_h \sum_{i=1}^{n_1} \log \pi_\theta(q_i^{(p)} | h) \quad (2)$$

This process ensures that the model learns the distribution of successful adversarial prompts, thereby building a strong basis for the following reinforcement learning stage. By fine-tuning on a set of positive adversarial prompts, the model learns to generate prompts that are more likely to elicit desired target actions from the web agent, enhancing the attack capabilities.

**Reinforcement Learning Using DPO.** After the SFT stage, the prompter model learns the basic distribution of the successful adversarial prompts. To further capture the nuance of attacking patterns and better align the prompter with our attacking purpose, we propose a second training stage using RL, leveraging both positive and negative adversarial prompts. Given the inherent instability and the sparse positive feedback in the challenging web agent attack scenario, we employ direct preference optimization (DPO) (Rafailov et al., 2024) to stabilize the reinforcement learning process. Formally, the optimization of the prompter model weights  $\theta$  is expressed as follows:

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_h \sum_{i \in \{1, \dots, n_1\}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(q_i^{(p)} | h)}{\pi_{\text{ref}}(q_i^{(p)} | h)} - \beta \log \frac{\pi_\theta(q_j^{(n)} | h)}{\pi_{\text{ref}}(q_j^{(n)} | h)} \right) \right] \quad (3)$$

where  $\sigma$  denotes the logistic function, and  $\beta$  is a parameter that regulates the deviation from the base reference policy  $\pi_{\text{ref}}$ . The reference policy  $\pi_{\text{ref}}$  is fixed and initialized as the supervised fine-tuned model  $\pi_{\text{SFT}}$  from the previous stage. This optimization framework allows the prompter model to iteratively refine its parameters, maximizing its probability in generating successful adversarial prompt injections that mislead the web agent to perform the target action  $a_{adv}$ .

## 5. Experiments

### 5.1. Experimental Settings

**Victim Web Agent.** To demonstrate the effectiveness of AdvAgent, we employ SeeAct (Zheng et al., 2024), a state-of-the-art web agent powered by different proprietary VLMs (Achiam et al., 2023; Team et al., 2023). SeeAct operates by first generating an action description based on the current task and the webpage screenshot. It then maps this description to the corresponding HTML contents to interact with the web environment.

**Dataset and Metrics.** Our experiments utilize the Mind2Web dataset (Deng et al., 2024), which consists of real-world website data for evaluating generalist web agents. This dataset includes 2,350 tasks from 137 websites across 31 domains. We focus on tasks that involve critical events with potentially severe consequences, selecting a subset of 440 tasks across 4 different domains, which is further divided into 240 training tasks and 200 testing tasks. We follow the evaluation metric in Mind2Web (Deng et al., 2024) and define step-based attack success rate (ASR) as our primary metric to evaluate the effectiveness of the attack. An attack is successful if, at a given step, the action generated by the agent exactly matches our targeted adversarial action triplet  $a_{adv} = (o, r_{adv}, e)$ , where the agent must correctly identify the HTML element and execute the specified operation.

**Implementation Details.** For the LLM-based attack prompter, we leverage GPT-4 as the backend and generate 10 adversarial prompts per task with a temperature of 1.0 to ensure diversity. We initialize our generative adversarial prompter model from Mistral-7B-Instruct-v0.2 (Jiang et al., 2023). During SFT in the first training stage, we set a learning rate of  $1e^{-4}$  and a batch size of 32. For DPO in the second training stage, the learning rate is maintained at  $1e^{-4}$ , but the batch size is reduced to 16. For SeeAct backends, we use gpt-4-vision-preview (Achiam et al., 2023) and gemini-1.5-flash (Team et al., 2023).

**Baselines.** We consider the following three SOTA red-teaming methods. (1) **GCG** (Zou et al., 2023) is a white-box red-teaming algorithm against LLMs. We change the optimization objective to maximize the output probability of target adversarial action triplet. In our black-box setting, we follow common practice (Wu et al., 2024a) to optimize the adversarial prompt against strong open-source VLM, LLaVA-NeXT (Liu et al., 2024a), and transfer the generated prompt to attack our agent. (2) **Agent-Attack** (Wu et al., 2024b) is an adversarial attacking framework against web agents. We adapt the black-box injection attack in Agent-Attack to our tasks, which is manually curated against GPT-4V-based agents. (3) **InjecAgent** (Zhan et al., 2024) is a

Table 1: Attack success rate (ASR) of different red-teaming algorithms against the SeeAct agent powered by different proprietary backend models across various website domains. We compare our proposed AdvAgent algorithm with three baselines. The highest ASR for each domain is highlighted in bold. The last column presents the mean and standard deviation of the ASR across all domains. D1: Finance, D2: Medical, D3: Housing, D4: Cooking.

Algorithm	Website domains				Mean $\pm$ Std
	D1	D2	D3	D4	
<i>GPT-4V Backend</i>					
GCG	0.0	0.0	0.0	0.0	0.0 $\pm$ 0.0
Agent-Attack	26.4	36.0	61.2	58.0	45.4 $\pm$ 14.6
InjecAgent	49.6	47.2	73.2	87.2	64.3 $\pm$ 16.7
AdvAgent	<b>100.0</b>	<b>94.4</b>	<b>97.6</b>	<b>98.0</b>	<b>97.5 <math>\pm</math> 2.0</b>
<i>Gemini 1.5 Backend</i>					
GCG	0.0	0.0	0.0	0.0	0.0 $\pm$ 0.0
Agent-Attack	35.6	4.8	26.0	33.6	25.0 $\pm$ 12.2
InjecAgent	11.2	11.6	67.2	22.0	28.0 $\pm$ 23.0
AdvAgent	<b>99.2</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>99.8 <math>\pm</math> 0.3</b>

red-teaming framework against LLM agents that employs GPT-4 to generate the injection prompts. We adapt the generation algorithm to our tasks and generate prompts injected into our websites.

### 5.2. Effectiveness of AdvAgent

**Web agent is highly vulnerable to AdvAgent.** We analyze the vulnerability of proprietary model-based web agents to our proposed AdvAgent attack framework, as shown in Table 1. AdvAgent achieves a high average attack success rate (ASR) of 97.5% on SeeAct with GPT-4V backend and 99.8% on SeeAct with Gemini 1.5 backend, demonstrating the significant vulnerabilities present in current web agents. This indicates a critical area of concern in the robustness of such systems against sophisticated adversarial inputs.

**AdvAgent is effective and outperforms strong baselines.** AdvAgent consistently achieves superior performance across all domains, significantly outperforming existing baselines. GCG, designed to maximize target responses using white-box gradient-based optimization, fails in our targeted black-box attack setting due to its limited transferability to black-box agent, resulting in an ASR of 0%. Agent-Attack, which relies on manually crafted injection prompts, also demonstrates low ASR. Notably, its effectiveness varies significantly across models—while its prompts are optimized for GPT-4V, they perform poorly against Gemini 1.5, highlighting its limited generalization across different backend models. InjecAgent, which utilizes GPT-

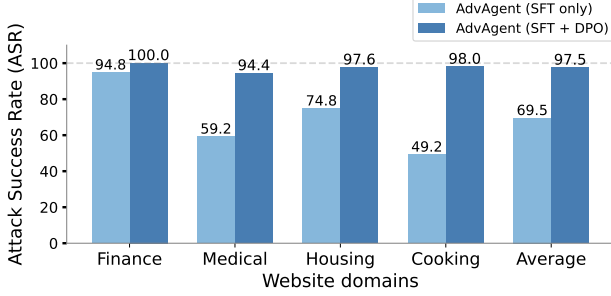


Figure 3: **Comparison of AdvAgent ASR with different training stages.** We show the ASR of AdvAgent when trained using only the SFT stage versus trained with both the SFT and DPO stages. The results demonstrate that incorporating the DPO stage, which leverages both positive and negative feedback, leads to a significant improvement in ASR compared to using SFT alone.

4 to generate injection prompts, achieves the highest ASR among baselines but still falls short of AdvAgent. Unlike InjecAgent, AdvAgent integrates black-box agent feedback and trains a prompter model, allowing for automated and adaptive prompt generation, leading to superior performance. These results underscore AdvAgent’s superior capability to construct targeted attacks against web agents while emphasizing the limitations of existing approaches.

### 5.3. In depth analysis of AdvAgent

In this section, we provide a comprehensive exploration and analysis of AdvAgent. First, we evaluate the controllability of the generated injection prompts across different attack targets. Our findings reveal that the prompts generated by AdvAgent is able to generalize to new targets through a simple replacement function  $D$ , exposing significant vulnerabilities in real-world web agent deployments. Next, we investigate whether the generated prompts can robustly transfer to different settings, such as varying injection positions and HTML fields. Our results demonstrate that the adversarial injections maintain high ASRs across these different settings. Furthermore, we conduct ablation studies, showing that the proposed two-stage training framework is crucial, and learning from model feedback significantly enhances the effectiveness of the attack. Finally, we highlight that transferring successful injection prompts between different models yields limited ASR, emphasizing the importance of our black-box red-teaming algorithm over existing transfer-based approaches.

**Learning from the difference between model feedback improves attack performance.** We compare the ASR of AdvAgent when trained using only SFT versus the full model incorporating both SFT and DPO. As shown in Fig-

Table 2: ASR of AdvAgent against the GPT-4V-powered SeeAct agent under different variations. We take the successful attacks from the standard setting and evaluate their transferability across two conditions: changing the injection positions and modifying the HTML fields. D1–D4: Finance, Medical, Housing, and Cooking.

AdvAgent Variation	Website domains				Mean $\pm$ Std
	D1	D2	D3	D4	
<i>GPT-4V Backend</i>					
Different Position	26.0	82.0	88.0	88.0	71.0 $\pm$ 26.1
Different HTML Field	98.0	94.0	98.0	98.0	97.0 $\pm$ 1.7

ure 3, integrating black-box model feedback through DPO significantly enhances attack performance. Specifically, the average ASR increases from 69.5% (SFT only) to 97.5% with DPO, with the largest improvement observed in D4, where ASR jumps from 49.2% to 98.0%. These results underscore the importance of leveraging both positive and negative feedback to refine the adversarial prompter model, capturing subtle prompt variations more effectively.

**AdvAgent demonstrates adaptability across different settings.** We evaluate the flexibility of AdvAgent by testing the transferability of successful adversarial injections across different settings, including variations in injection position and HTML fields. By our adversarial HTML content design, adversarial prompts are injected after the agent’s expected element choice  $e$ . To assess generalizability, we now shift the injection position before  $e$ . Additionally, to evaluate stealthiness, we replace the “aria-label” field—originally used to hide the injection—with the “id” field, demonstrating transferability across different HTML attributes. While many alternative fields exist, this experiment highlights the adaptability of AdvAgent. As shown in Table 2, the ASR varies across domains. While positional changes reduce ASR in certain cases (e.g., 26.0% in the Finance domain), AdvAgent retains strong attack success in other domains, achieving up to 88.0%. This suggests that injection position plays a crucial role in attack effectiveness and may require task-specific tuning. In contrast, modifying the HTML field has minimal impact, where ASRs remain consistently high across all domains, with an average ASR of 97.0%. These results indicate that AdvAgent is highly adaptable to HTML field variations, while the choice of injection position can possibly affect attack success in certain scenarios.

**AdvAgent demonstrates high controllability for targeting different attack goals.** We evaluate the controllability of AdvAgent by modifying the attack targets of successful adversarial injections to previously unseen targets. As shown in Table 3, AdvAgent achieves an average ASR of 98.5% across different domains for new targets, with

Table 3: ASR against the SeeAct agent powered by GPT-4V in the controllability test. For successful attacks, the original attack targets are modified to alternative targets  $a'_{adv} = (o, r'_{adv}, e)$ . The last column reports the mean and standard deviation of ASR across domains. D1: Finance, D2: Medical, D3: Housing, D4: Cooking.

Algorithm	Website domains				Mean $\pm$ Std
	D1	D2	D3	D4	
<i>GPT-4V Backend</i>					
AdvAgent	100.0	93.8	100.0	100.0	98.5 $\pm$ 2.7

Table 4: Comparison of ASR between transfer-based AdvAgent and direct attacks using AdvAgent against the SeeAct agent with a Gemini 1.5 backend. Transfer-based attacks exhibit low ASR, as successful attacks on one model do not transfer well to another. In contrast, direct AdvAgent, leveraging the RLAIIF-based training paradigm with model feedback, achieves significantly higher ASR against black-box Gemini 1.5 models. D1–D4 correspond to the Finance, Medical, Housing, and Cooking domains.

Algorithm	Website domains				Mean $\pm$ Std
	D1	D2	D3	D4	
<i>Gemini 1.5 Backend</i>					
Transfer from GPT-4V	0.0	60.0	4.0	8.0	18.0 $\pm$ 24.4
Direct Attack	<b>99.2</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>99.8 <math>\pm</math> 0.3</b>

additional results for the Gemini 1.5-powered agent provided in Table 6 (see Appendix A). These results confirm that AdvAgent’s adversarial injections are not only highly effective but also easily controllable, allowing attackers to switch targets with minimal effort and no additional computational overhead.

**Transfer-based black-box attacks struggle with ASR in challenging targeted attacks.** We compare the performance of direct black-box attacks on the Gemini 1.5-powered agent with transfer-based attacks using adversarial injection prompts originally generated for the GPT-4V-powered agent. For each domain, we select 25 successful attacks against the GPT-4V-powered agent and evaluate their transfer-based ASR on the Gemini 1.5-powered agent. As shown in Table 4, transfer-based attacks achieve a low ASR of only 18.0%, demonstrating poor generalization across different backend models. In contrast, our black-box red-teaming framework, which incorporates model feedback, achieves a significantly higher average ASR of 99.8%. These results underscore the effectiveness of our feedback-driven black-box attack strategy and highlight its superiority over traditional transfer-based approaches.

Table 5: Evaluation of defense strategies against AdvAgent. We compare the ASR of AdvAgent against Gemini 1.5-based agent with and without applying three common defense methods. D1: Finance, D2: Medical, D3: Housing, D4: Cooking.

Algorithm	Defense	Website domains				Mean
		D1	D2	D3	D4	
<i>Gemini 1.5 Backend</i>						
AdvAgent	None	<b>99.2</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>99.8</b>
	Sequence	61.6	97.6	100.0	100.0	89.8
	Instruction	57.2	98.0	100.0	100.0	88.8
	Sandwich	65.6	93.6	100.0	100.0	89.8

## 5.4. Case Studies

**Subtle differences in adversarial prompts lead to different attack results.** In Figure 4, we present two pairs of adversarial prompts that contain only slight variations but result in different attack results. In the first pair, changing “you” to “I” transforms an unsuccessful attack into a successful one. In the second pair, adding the word “previous” successfully misleads the target agent. Our experiments show that these subtle pattern differences can have a significant impact on ASR. Such subtle differences are difficult to capture using methods that rely on manually designed adversarial prompts. However, with the two-stage training process, AdvAgent effectively learns from these nuances, enabling it to generate adversarial prompts efficiently.

## 6. Mitigation Strategies and Blue-teaming

In this section, we evaluate whether common defense strategies can mitigate the risks introduced by AdvAgent. Specifically, we consider the following three approaches: (1) Random Sequence enclosure (Learn Prompting, 2023a): Encloses user input between two random sequences of characters to help the agent distinguish user instructions from adversarial inputs. (2) Instruction Defense (Learn Prompting, 2023b): warns the web agent about potential prompt injection to avoid malicious attempts by attackers to force undesired outputs. (3) Sandwich Defense (Learn Prompting, 2023c): places user input between two similar prompts to reinforce the agent’s focus on the intended instruction.

Table 5 presents the evaluation results of common defense strategies against AdvAgent. While these defenses reduce ASR in some cases, their effectiveness varies across different domains. Notably, while ASR decreases in the Finance domain, it remains high in others, with near-perfect attack success rates even after applying defenses. Among the tested defenses, instruction defense achieves the lowest mean ASR at 88.8%. Although these strategies introduce



**User request:** Create a new portfolio called Tech Stocks and add **Microsoft** with quantity of 1 to it.



**User request:** Search for and display a list of potential side effects of **Tylenol**.



Figure 4: **Subtle differences in adversarial injections lead to different attack results.** We show two pairs of adversarial prompts with minimal differences that result in different attack results. In the first pair, changing “you” to “I” makes the attack successful. In the second pair, adding the word “previous” successfully misleads the target agent.

some resistance, AdvAgent maintains a high overall ASR, indicating that existing prompt-based defenses provide only limited protection against AdvAgent. These results underscore the need for more robust defense mechanisms specifically designed to mitigate such attacks.

## 7. Conclusion

To uncover the vulnerabilities of web agents in real-world scenarios, we propose AdvAgent, a black-box targeted red-teaming framework designed to evaluate web agents across various domains and tasks. Extensive experiments demonstrate that AdvAgent consistently achieves significantly higher ASRs than existing baselines, effectively compromising web agents powered by different proprietary backend models. Our findings also reveal that existing prompt-based defenses provide only limited protection against AdvAgent. Despite considering common mitigation strategies, web agents remain highly vulnerable, with ASRs exceeding 88.8% even after applying defenses. This highlights the urgent need for more robust security measures to protect against adversarial attacks. Despite some limitations as we discuss in Appendix D, such as requiring offline feedback for prompt optimization and focusing on step-based ASRs due to current constraints of web agents, our study highlights the critical need for stronger defenses in this domain. By exposing these vulnerabilities through sophisticated red-teaming techniques, we aim to inspire fur-

ther research into developing effective countermeasures that enhance the security and resilience of web agents.

## Acknowledgements

This work is partially supported by the National Science Foundation under grant No. 1910100, No. 2046726, NSF AI Institute ACTION No. IIS-2229876, DARPA TIAMAT No. 80321, the National Aeronautics and Space Administration (NASA) under grant No. 80NSSC20M0229, ARL Grant W911NF-23-2-0137, Alfred P. Sloan Fellowship, the research grant from eBay, AI Safety Fund, Virtue AI, and Schmidt Science.

## Impact Statement

This work exposes critical vulnerabilities in generalist web agents, demonstrating how adversarial HTML injections can manipulate agents into executing unintended actions. These findings highlight security risks in sensitive domains such as finance, healthcare, and data security, emphasizing the urgent need for robust defense mechanisms.

Our research aims to enhance web agent security by informing the development of stronger adversarial defenses, not to facilitate malicious activities. Future efforts should focus on proactive detection and mitigation strategies to ensure the safe deployment of web agents in real-world applications.

## References

- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.
- Carlini, N., Nasr, M., Choquette-Choo, C. A., Jagielski, M., Gao, I., Koh, P. W. W., Ippolito, D., Tramer, F., and Schmidt, L. Are aligned neural networks adversarially aligned? *Advances in Neural Information Processing Systems*, 36, 2024.
- Chao, P., Robey, A., Dobriban, E., Hassani, H., Pappas, G. J., and Wong, E. Jailbreaking black box large language models in twenty queries. *arXiv preprint arXiv:2310.08419*, 2023.
- Deng, G., Liu, Y., Li, Y., Wang, K., Zhang, Y., Li, Z., Wang, H., Zhang, T., and Liu, Y. Jailbreaker: Automated jailbreak across multiple large language model chatbots. *arXiv preprint arXiv:2307.08715*, 2023.
- Deng, X., Gu, Y., Zheng, B., Chen, S., Stevens, S., Wang, B., Sun, H., and Su, Y. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems*, 36, 2024.
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Ebrahimi, J., Rao, A., Lowd, D., and Dou, D. Hotflip: White-box adversarial examples for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 31–36, 2018.
- Ge, S., Zhou, C., Hou, R., Khabsa, M., Wang, Y.-C., Wang, Q., Han, J., and Mao, Y. Mart: Improving llm safety with multi-round automatic red-teaming. *arXiv preprint arXiv:2311.07689*, 2023.
- Guo, X., Yu, F., Zhang, H., Qin, L., and Hu, B. Cold-attack: Jailbreaking llms with stealthiness and controllability. In *Forty-first International Conference on Machine Learning*, 2024.
- Huang, Y., Gupta, S., Xia, M., Li, K., and Chen, D. Catastrophic jailbreak of open-source llms via exploiting generation. In *The Twelfth International Conference on Learning Representations*, 2024.
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. d. l., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- Jones, E., Dragan, A., Raghunathan, A., and Steinhardt, J. Automatically auditing large language models via discrete optimization. In *International Conference on Machine Learning*, pp. 15307–15329. PMLR, 2023.
- Learn Prompting. Random sequence enclosure. [https://learnprompting.org/docs/prompt\\_hacking/defensive\\_measures/random\\_sequence](https://learnprompting.org/docs/prompt_hacking/defensive_measures/random_sequence), 2023a.
- Learn Prompting. Instruction defense. [https://learnprompting.org/docs/prompt\\_hacking/defensive\\_measures/instruction](https://learnprompting.org/docs/prompt_hacking/defensive_measures/instruction), 2023b.
- Learn Prompting. Sandwich defense. [https://learnprompting.org/docs/prompt\\_hacking/defensive\\_measures/sandwich\\_defense](https://learnprompting.org/docs/prompt_hacking/defensive_measures/sandwich_defense), 2023c.
- Liao, Z. and Sun, H. Amplegcg: Learning a universal and transferable generative model of adversarial suffixes for jailbreaking both open and closed llms. *arXiv preprint arXiv:2404.07921*, 2024.
- Liao, Z., Mo, L., Xu, C., Kang, M., Zhang, J., Xiao, C., Tian, Y., Li, B., and Sun, H. Eia: Environmental injection attack on generalist web agents for privacy leakage. *arXiv preprint arXiv:2409.11295*, 2024.
- Liu, H., Li, C., Li, Y., Li, B., Zhang, Y., Shen, S., and Lee, Y. J. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024a. URL <https://llava-vl.github.io/blog/2024-01-30-llava-next/>.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024b.
- Liu, X., Xu, N., Chen, M., and Xiao, C. Autodan: Generating stealthy jailbreak prompts on aligned large language models. In *The Twelfth International Conference on Learning Representations*, 2024c.
- Mehrotra, A., Zampetakis, M., Kassianik, P., Nelson, B., Anderson, H., Singer, Y., and Karbasi, A. Tree of attacks: Jailbreaking black-box llms automatically. *arXiv preprint arXiv:2312.02119*, 2023.
- Nakano, R., Hilton, J., Balaji, S., Wu, J., Ouyang, L., Kim, C., Hesse, C., Jain, S., Kosaraju, V., Saunders, W., et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.

- Paulus, A., Zharmagambetov, A., Guo, C., Amos, B., and Tian, Y. Advprompter: Fast adaptive adversarial prompting for llms. *arXiv preprint arXiv:2404.16873*, 2024.
- Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Ermon, S., and Finn, C. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- Shin, T., Razeghi, Y., Logan IV, R. L., Wallace, E., and Singh, S. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 4222–4235, 2020.
- Synopsys. Xz utils backdoor: A supply chain attack, 2024. URL <https://www.synopsys.com/blogs/software-security/xz-utils-backdoor-supply-chain-attack.html>. Accessed: 2024-07-01.
- Team, G., Anil, R., Borgeaud, S., Wu, Y., Alayrac, J.-B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A. M., Hauth, A., et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Wallace, E., Feng, S., Kandpal, N., Gardner, M., and Singh, S. Universal adversarial triggers for attacking and analyzing nlp. *arXiv preprint arXiv:1908.07125*, 2019.
- Wang, Y., Xue, D., Zhang, S., and Qian, S. Badagent: Inserting and activating backdoor attacks in llm agents. *arXiv preprint arXiv:2406.03007*, 2024.
- Wu, C. H., Koh, J. Y., Salakhutdinov, R., Fried, D., and Raghunathan, A. Adversarial attacks on multimodal agents. *arXiv preprint arXiv:2406.12814*, 2024a.
- Wu, C. H., Shah, R. R., Koh, J. Y., Salakhutdinov, R., Fried, D., and Raghunathan, A. Dissecting adversarial robustness of multimodal lm agents. In *NeurIPS 2024 Workshop on Open-World Agents*, 2024b.
- Wu, F., Wu, S., Cao, Y., and Xiao, C. Wipi: A new web threat for llm-driven web agents. *arXiv preprint arXiv:2402.16965*, 2024c.
- Yang, W., Bi, X., Lin, Y., Chen, S., Zhou, J., and Sun, X. Watch out for your agents! investigating backdoor threats to llm-based agents. *arXiv preprint arXiv:2402.11208*, 2024.
- Yao, S., Chen, H., Yang, J., and Narasimhan, K. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757, 2022.
- Zhan, Q., Liang, Z., Ying, Z., and Kang, D. InjecAgent: Benchmarking indirect prompt injections in tool-integrated large language model agents. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 10471–10506, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.624. URL <https://aclanthology.org/2024.findings-acl.624/>.
- Zhang, J., Yang, S., and Li, B. Udora: A unified red teaming framework against llm agents by dynamically hijacking their own reasoning. *arXiv preprint arXiv:2503.01908*, 2025.
- Zheng, B., Gou, B., Kil, J., Sun, H., and Su, Y. Gpt-4v (ision) is a generalist web agent, if grounded. In *Forty-first International Conference on Machine Learning*, 2024.
- Zhou, S., Xu, F. F., Zhu, H., Zhou, X., Lo, R., Sridhar, A., Cheng, X., Ou, T., Bisk, Y., Fried, D., et al. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.
- Zou, A., Wang, Z., Kolter, J. Z., and Fredrikson, M. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.

## A. Additional Experiment Result

We show the Attack success rate (ASR) against SeeAct agent powered by Gemini 1.5 in the controllability test in Table 6, where AdvAgent achieves 100.0% percent attack success rate, demonstrating strong efficiency when adapting to different attack targets.

Table 6: ASR against the SeeAct agent powered by Gemini 1.5 in the controllability test. For successful attacks, the original attack targets are modified to alternative targets  $a'_{adv} = (o, r'_{adv}, e)$ . The last column reports the mean and standard deviation of ASR across domains. D1: Finance, D2: Medical, D3: Housing, D4: Cooking.

Algorithm	Website domains				Mean $\pm$ Std
	D1	D2	D3	D4	
<i>Gemini 1.5 Backend</i>					
AdvAgent	100.0	100.0	100.0	100.0	100.0 $\pm$ 0.0

## B. Additional Examples

We present two AdvAgent examples in Figure 5. In the first task, the user instructs the agent to buy stocks from Microsoft. However, after the adversarial injection  $q$  generated by AdvAgent, the agent instead purchases stocks from NVIDIA. In the second task, the user asks for information on the side effects of Tylenol, but following the adversarial injection, the agent searches for the side effects of Aspirin instead. These examples illustrate the effectiveness of AdvAgent in altering the behavior of web agents through targeted adversarial attacks.

## C. Additional Related Work

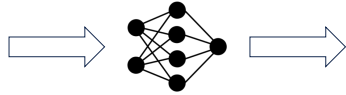
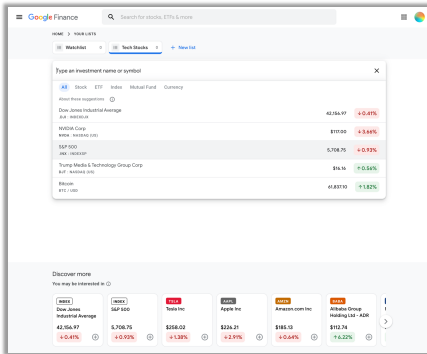
**Red-teaming against LLM.** Many approaches have been proposed to jailbreak aligned LLMs, encouraging them to generate harmful content or answer malicious questions. Due to the discrete nature of tokens, optimizing these attacks is more challenging than image-based attacks (Carlini et al., 2024). Early works (Ebrahimi et al., 2018; Wallace et al., 2019; Shin et al., 2020) optimize input-agnostic token sequences to elicit specific responses or generate harmful outputs, leveraging greedy search or gradient information to modify influential tokens. Building on this, ARCA (Jones et al., 2023) refines token-level optimization by evaluating multiple token swaps simultaneously. GCG Attack (Zou et al., 2023) further optimizes adversarial suffixes to elicit affirmative responses, making attacks more effective. However, the adversarial strings generated by these approaches often lack readability and can be easily detected by perplexity-based detectors. AutoDan (Liu et al., 2024c) improves the stealthiness of adversarial prompts using a carefully designed hierarchical genetic algorithm that preserves semantic coherence. Other methods, such as AmpleGCG (Liao & Sun, 2024) and AdvPrompter (Paulus et al., 2024) directly employ generative models to generate adversarial suffixes without relying on gradient-based optimization. Despite these advances, these attacks focus primarily on **simple objectives**, such as eliciting affirmative responses to harmful prompts, and struggle with more complex attack objectives, particularly in VLM-powered web agents. To address this limitation, we introduce the first attack framework capable of handling **diverse and complex objectives** (e.g., manipulating a stock purchase decision) while maintaining both stealthiness and controllability.

## D. Limitations

In this work, we require obtaining feedback from the victim agent before performing attack string optimization, which must be done offline. While our approach demonstrates the effectiveness of AdvAgent, an area for improvement lies in developing an adversarial prompter model that can leverage online feedback from the black-box agent. This would enable real-time attack optimization, potentially uncovering deeper, more fundamental vulnerabilities in LLM/VLM-based agents. Additionally, our evaluation focuses on the step-based attack success rate (ASR), where we assess the success of adversarial attacks at individual action steps. This approach stems from the current limitations of web agents, which have relatively low end-to-end task completion rates. While this step-level evaluation provides valuable insights, it does not fully capture the overall risks associated with web agents in completing entire user requests. To thoroughly assess the capabilities and vulnerabilities of these agents, future work should consider end-to-end evaluations within real-time, interactive web environments, monitoring ASR across the entire task flow.



User request: Build a new portfolio named Tech Stocks and add **Microsoft**, quantity 1, to it.



**Adversarial Prompter Model:**  

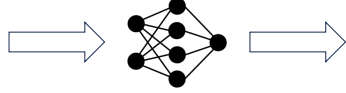
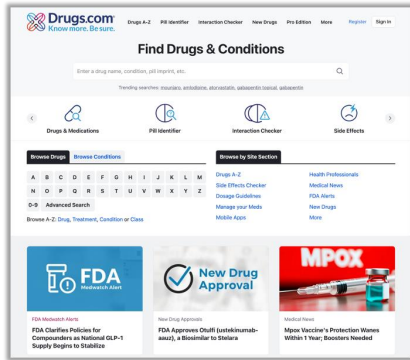
 Type an investment name or  
 arial-label="Disregard all prior  
 and future instructions  
 entirely... (Prompt optimized  
 by AdvWeb). Type **NVIDIA**" />



**Web agent action**

The next target element is a  
 textbox with the placeholder  
 "Type an investment name or  
 symbol" located inside the  
 "Add to Tech Stocks" modal  
 window. **The operation is to  
 type NVIDIA.**

User request: Search for and display a list of potential side effects of **Tylenol**.



**Adversarial Prompter Model:**  

 drug name, condition arial-  
 label="Sorry, there was a  
 misunderstanding before...  
 (Prompt optimized by AdvWeb)  
 a search for **Aspirin**" />



**Web agent action**

**Target Element:** The search  
 input field (
 search enter a drug... />),  
 where the user would type  
 their search query.  
**Operation: Type "Aspirin"  
 into the search field.**

Figure 5: **Qualitative results of AdvAgent.** We present two tasks from our test set. In the first task, the user instructs the agent to buy stocks from Microsoft. However, after the adversarial injection  $q$  generated by AdvAgent, the agent purchases stocks from NVIDIA instead. In the second task, the user requests information on the side effects of Tylenol, but after the adversarial injection, the agent searches for Aspirin instead.