

# Peering Behind the Shield: Guardrail Identification in Large Language Models

Ziqing Yang<sup>1</sup>, Yixin Wu<sup>1</sup>, Rui Wen<sup>2</sup>, Michael Backes<sup>1</sup>, Yang Zhang<sup>1\*</sup>

<sup>1</sup>CISPA Helmholtz Center for Information Security, <sup>2</sup>Institute of Science Tokyo,  
{ziqing.yang, yixin.wu}@cispa.de, wen@comp.isct.ac.jp, {director, zhang}@cispa.de

## Abstract

With the rapid adoption of large language models (LLMs), conversational AI agents have become widely deployed across real-world applications. To enhance safety, these agents are often equipped with guardrails that moderate harmful content. Identifying the guardrails in an agent thus becomes critical for adversaries to understand the system and design guard-specific attacks. In this work, we introduce **AP-Test**, a novel approach that leverages guard-specific adversarial prompts to detect the identity of guardrails deployed in black-box AI agents. Our method addresses key challenges in this task, including the influence of safety-aligned LLMs and other guardrails, as well as a lack of principled decision-making strategies. **AP-Test** employs two complementary testing strategies, input and output guard tests, and a new metric, match score, to enable robust identification. Experiments across diverse agents and four open-source guardrails demonstrate that **AP-Test** achieves perfect classification accuracy in multiple scenarios. Ablation studies further highlight the necessity of our proposed components. Our findings reveal a practical path toward guardrail identification in real-world AI systems.<sup>1</sup>

## 1 Introduction

Human-AI conversations have been significantly advanced by the rapid development of large language models (LLMs). These conversational AI agents are now extensively deployed across various domains, including customer service (DeepSeek-AI et al., 2025; Anthropic; OpenAI, 2024), education (Neumann et al., 2024), and healthcare (Peng et al., 2023). Such widespread use also raises severe security concerns, such as jailbreak attacks (Shen et al., 2024; Zou et al., 2023; Liu

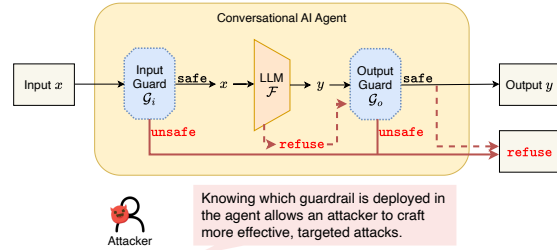


Figure 1: A conversational AI agent equipped with guardrails. Normally, the agent responds to the input following the black line. Once the guardrail flags the input or generated response unsafe, the agent would refuse the query (solid red line). The dashed red line signifies the LLM’s refusal to the input per its internal safety alignment.

et al., 2023; Zhao et al., 2024) and prompt injection attacks (Liu et al., 2024a; Zhan et al., 2024; Debenedetti et al., 2024).

Safety guardrails (Inan et al., 2023; Zeng et al., 2024) are thus developed to further moderate the input/output content of the AI agents, as shown in Figure 1. They are often fine-tuned on top of LLMs using annotated datasets that cover a broad spectrum of safety risks (Inan et al., 2023; Han et al., 2024; Zeng et al., 2024). In addition, a growing number of high-quality open-source guardrails, such as LlamaGuard (Inan et al., 2023), WildGuard (Han et al., 2024), and ShieldGemma (Zeng et al., 2024), have become widely available and are increasingly integrated into AI Agents. This largely challenges attackers to bypass the safety mechanisms in AI agents; thus, it is crucial to identify the guardrails deployed in the AI agent. Once a guardrail is identified, it can be treated as a white-box component, facilitating guard-specific attacks that are shown to outperform black-box methods (Mangaokar et al., 2024; Zhang et al., 2024a).

Considering that guardrails are often built on top of LLMs, a natural question arises: can we simply apply LLM identification methods (Zeng et al.,

\*Corresponding author.

<sup>1</sup>Our code of **AP-Test** is available at <https://github.com/TrustAILab/AP-Test>.

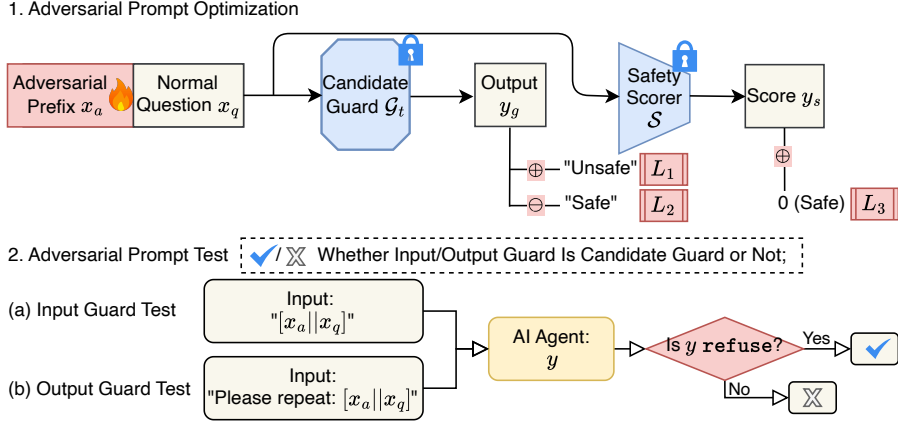


Figure 2: Framework of our **AP-Test**. We first perform *adversarial prompt optimization* based on the candidate guardrail with our tailored loss function. Then, we conduct *adversarial prompt tests* by querying the AI agent with the adversarial prompts to determine whether the candidate guardrail exists in the agent.

2023; Gubri et al., 2024; Jin et al., 2024; Yang et al., 2025) to detect the presence of specific guardrails in an AI agent? LLM identification aims to determine the origin or identity of an LLM, primarily to address concerns about unauthorized use. The common approach relies on optimizing adversarial prompts that elicit distinctive responses from the candidate LLM (Gubri et al., 2024; Jin et al., 2024). However, this line of work falls short due to several challenges in guardrail identification: **C1** The output space of guardrails is extremely limited (e.g., binary safe/unsafe responses), so optimizing input triggers toward a target response is much easier. This makes it hard to be specific to one guardrail. **C2** AI agents usually allow query-only access and do not provide knowledge of the underlying safety mechanisms (OpenAI, 2024; Anthropic; Github; DeepSeek-AI et al., 2025), revealing little information about the guardrail. Mechanisms such as additional guardrails or safety-aligned LLMs may also obscure the existence of the candidate guardrail. **C3** Guardrails are integrated into AI agents with unknown deployment stages; when applied at the output stage, their input depends on the LLM’s outputs. In this sense, it is challenging to input the adversarial prompts into the output guardrail. **C4** Existing LLM identification methods lack principled strategies for selecting identification thresholds, limiting their robustness and reliability in practical settings.

By solving the above challenges, we propose **AP-Test**, which utilizes **A**dversarial **P**rompts to **T**est the identity of the safety guardrail deployed in a black-box AI agent. As shown in Figure 2, our ap-

proach begins by probing the AI agent with guard-specific adversarial prompts, which are designed to be flagged as unsafe by a specific candidate guardrail while remaining safe according to others. Such optimization exploits the limited output space of guardrails (**C1**) and mitigates the influence of safety-aligned LLMs and other safety guardrails (**C2**) via a tailored loss function. We then perform identification using two complementary strategies: the *input guard test* and the *output guard test*. Together, these tests cover scenarios where guardrails are deployed at either the input or output stage of an agent, thereby addressing **C3**. Intuitively, the output guard test asks the agent to repeat the adversarial prompt, so that the output guard will receive the LLM’s response that contains the repeated text. To make identification without requiring additional guardrails (**C4**), we design a novel metric, match score, to measure how the safety mechanism in the agent matches the candidate guardrail. A larger match score indicates a stronger likelihood of the candidate guardrail being deployed.

To show the practicability of our **AP-Test**, we conduct extensive experiments on four different candidate guardrails under diverse scenarios, covering two popular LLMs and 10 guardrails. Results demonstrate that the proposed attack method accurately identifies guardrails in various AI agents for both input and output guardrails. Specifically, the WildGuard-specific (Han et al., 2024) input guard test on both Llama3.1-based (Dubey et al., 2024) and GPT4o-based (OpenAI, 2024) agents achieves a perfect classification accuracy, i.e.,  $Acc = 1.00$ . That is, our **AP-Test** successfully identifies the exist-

tence of WildGuard in all evaluated agents, indicating its effectiveness. Our **AP-Test** also successfully probes the candidate guardrail in more complex agents, e.g., containing two different guardrails. This showcases the practicality of our method.

Moreover, to analyze the impact of our proposed method, we perform an ablation study to examine the role of each component, particularly our loss terms designed to optimize adversarial queries and the existence of the query set. Our findings reveal that our loss terms and the query set are crucial for the identification. For example, without the loss ensuring the prompt remains safe for other guardrails ( $L_3$ ), **AP-Test** tends to misidentify that the LlamaGuard3 (Chi et al., 2024) is used in the agent that is only equipped with WildGuard as the input guardrail.

Overall, our contributions are as follows:

- We propose **AP-Test**, the first method that uses guard-specific adversarial prompts to identify guardrails in black-box AI agents, overcoming core challenges under this setting.
- Experiments show the effectiveness and robustness of **AP-Test** on four candidate guardrails on various agents under diverse scenarios.
- Our ablation study demonstrates the importance of each component in our proposed method, showing that their removal significantly degrades identifying performance.

## 2 Background and Related Work

**LLM Security Risks.** The rapid advancement of LLMs provides users with significant convenience but also raises critical security concerns (Zhao et al., 2024; Gong et al., 2023; Zhang et al., 2024b; Wen et al., 2024; Li et al., 2023). Among these concerns, jailbreak attacks (Shen et al., 2024; Zou et al., 2023; Liu et al., 2023; Zhao et al., 2024) pose a major threat by bypassing built-in safety mechanisms, enabling models to generate restricted or harmful content that violates usage policies (Google, 2024; Meta; OpenAI; Amazon, 2025). Previous studies analyze existing jailbreak strategies, particularly focusing on in-the-wild jailbreak prompts that are manually crafted in real-world scenarios (Shen et al., 2024). More recent studies introduce automated jailbreak generators, such as AutoDAN (Liu et al., 2023), GCG (Zou et al., 2023),

and TAP (Mehrotra et al., 2023), which optimize adversarial prompts to evade safety measures and maximize attack success rates.

**Internal Safety Alignment.** LLM safety alignment refers to the process of ensuring that LLMs generate outputs that are consistent with human values. Most popular LLMs, such as Llama3 (Dubey et al., 2024), Gemma (Mesnard et al., 2024), Mistral (Jiang et al., 2023), and ChatGPT (OpenAI, 2024), are safety-aligned, either through reinforcement learning with human feedback (RLHF) (Liu et al., 2024b; Ji et al., 2023), or learning from curated datasets that contain safety-related data (Jiang et al., 2023; Mesnard et al., 2024). However, these safety-aligned LLMs are still exposed to security risks such as jailbreak attacks (Shen et al., 2024; Zou et al., 2023; Liu et al., 2023; Zhao et al., 2024) and prompt injection attacks (Liu et al., 2024a; Zhan et al., 2024; Debenedetti et al., 2024; Yu et al., 2023).

**External Safety Guardrail.** Safety guardrails are designed for moderating the input/output content of the LLMs (Inan et al., 2023; Chi et al., 2024; Ghosh et al., 2024; Han et al., 2024; Zeng et al., 2024) and serve as an external complement to safety alignment. For example, built on Llama2-7B (Touvron et al., 2023), LlamaGuard (Inan et al., 2023) is fine-tuned on their constructed dataset based on a safety risk taxonomy encompassing a range of safety risks. Subsequent versions, LlamaGuard2 (Meta, 2024) and LlamaGuard3 (Chi et al., 2024), further expand the safety risk taxonomy and dataset, leveraging state-of-the-art LLMs for fine-tuning, thereby strengthening their safeguard capabilities. Similarly, Aegis (Ghosh et al., 2024), WildGuard (Han et al., 2024), and ShieldGemma (Zeng et al., 2024) follow a comparable approach. Specifically, the Aegis series is further instruction-tuned on LlamaGuard based on their own dataset, and ShieldGemma is built upon Gemma2 (Mesnard et al., 2024).

**LLM Fingerprinting and Watermarking.** LLM identification focuses on identifying the origin of an LLM (Zeng et al., 2023; McGovern et al., 2024; Gubri et al., 2024; Jin et al., 2024; Yang et al., 2025). A common approach involves crafting LLM-specific adversarial prompts to guide the candidate LLM to produce target responses, which are then used for identification (Gubri et al., 2024; Jin et al., 2024). As state-of-the-art guardrails are often built upon LLMs (Inan et al., 2023; Han et al.,

2024; Zeng et al., 2024), guardrail identification is similar. However, a key difference lies in the limited and often masked output space of guardrails (C1). This makes it hard to directly adapt the target responses for identification. Moreover, existing LLM identification techniques typically rely on empirically chosen thresholds by querying an auxiliary set of LLMs (Gubri et al., 2024; Jin et al., 2024), which fails to be solely based on the candidate model. This limitation highlights the need for more principled decision-making strategies, as addressed by our method (C4).

### 3 Problem Statement

**Preliminary.** At the core of these AI agents lies a safety-aligned LLM  $\mathcal{F}$  that drives their functionality, but still suffers from security risks (Shen et al., 2024; Zou et al., 2023; Liu et al., 2023, 2024a; Yu et al., 2023). To ensure the security and compliance of these AI agents, additional mechanisms known as input and output guardrails  $\mathcal{G}$  are often implemented. As shown in Figure 1, the input guardrail  $\mathcal{G}_i$  evaluates user inputs  $x$  to determine whether they should be forwarded to the LLM. If the input  $x$  is deemed high-risk, policy-violating, or jailbreak prompts, i.e.,  $\mathcal{G}_i(x) = \text{unsafe}$ , the agent then returns refuse without processing it further. Otherwise, the safety-aligned LLM would react to the input  $x$ , either generating responses or returning refuse if it perceives unsafe. However, with attacks like jailbreak attacks, even though the prompt input may be considered safe by both the input guard  $\mathcal{G}_i$  and the LLM  $\mathcal{F}$ , the LLM’s response could still contain harmful content. Therefore, in real-world applications, solely monitoring input may not be sufficient, which motivates the deployment of output guardrails. The output guardrail  $\mathcal{G}_o$  monitors the LLM’s response  $y$  to avoid policy violations. If a response is non-compliant, i.e.,  $\mathcal{G}_o(x) = \text{unsafe}$ , the agent would withhold the output and return refuse; otherwise, the agent would output the response  $y$ .

**Goal.** The goal of guardrail identification is to determine whether a guardrail or its derivative is deployed in an AI agent. The derivative refers to further customized versions of the guardrail, such as those that are fine-tuned or instruction-tuned.

**Capability.** We assume black-box access to the AI agent, with no knowledge of the underlying LLM or the presence of input/output guardrails. We focus on open-source guardrails due to their

widespread adoption. According to HuggingFace, LlamaGuard3 (Chi et al., 2024) received 335,571 downloads in March 2025, while Aegis-Permissive (Ghosh et al., 2024) was downloaded 1,021,359 times. One might argue that developers could create entirely proprietary guardrails. However, given the widespread adoption and accessibility of open-source guardrails, it is often more practical to build on existing ones through further fine-tuning. As such, we also consider derivatives of popular open-source guardrails to better cover the space of real-world deployments.

**Problem Formulation.** We define the guardrail identification task as follows: *Given an AI agent and a candidate guardrail  $\mathcal{G}_t$ , guardrail identification aims to identify whether the candidate guardrail  $\mathcal{G}_t$  or its derivative is deployed in the AI agent.* The identification task consists of two tests: (1) The input guard test audits whether the candidate guardrail  $\mathcal{G}_t$  is deployed at the input stage of the AI agent; (2) The output guard test evaluates whether the candidate guardrail is present at the output stage.

## 4 Method

In this work, we propose **AP-Test**, which leverages guard-specific **A**dversarial **P**rompts to **T**est the identity of the input/output guardrail deployed in an AI agent. As shown in Figure 2, our framework consists of two phases: *adversarial prompt optimization* and *adversarial prompt test*.

### 4.1 Adversarial Prompt Optimization

The goal of the optimization phase is to optimize adversarial prompts that the target guardrail  $\mathcal{G}_t$  erroneously flags as unsafe, while all other guardrails and safety-aligned LLMs correctly classify them as safe. An adversarial prompt is constructed by concatenating an optimized adversarial prefix  $x_a$  with a normal query  $x_q \in Q$ , where  $Q$  is a query set. The query is used as a starting point to prevent over-rejection by other guardrails, which is proven effective in our ablation study (Section 6). For each query  $x_q \in Q$ , we optimize an adversarial prefix  $x_a$  using three loss terms, addressing two specific aspects of the desired behavior.

**Candidate Guardrail Adversarial Losses ( $L_1$  and  $L_2$ ).** As guardrails’ output space is usually limited (C1), these loss terms are designed to mislead the candidate guardrail  $\mathcal{G}_t$  into classifying the adversarial prompt as unsafe. Specifically,  $L_1$  en-

courages the candidate guardrail to classify the adversarial prompt as unsafe, while  $L_2$  penalizes the candidate guardrail for classifying the adversarial prompt as safe:

$$\begin{aligned} L_1 &= \sigma(\mathcal{G}_t(x_a||x_q), \text{unsafe}), \\ L_2 &= -\sigma(\mathcal{G}_t(x_a||x_q), \text{safe}), \end{aligned} \quad (1)$$

where  $\sigma(\cdot, \cdot)$  represents the cross-entropy loss. Together,  $L_1$  and  $L_2$  ensure that  $x_a$  effectively triggers the refusal mechanism of the candidate guardrail. Conceptually, these two losses are identical, but our ablation study (Section 6) demonstrates that the synergy of these two losses slightly outperforms solely deploying a single one of them.

**Cross-Guardrail Compatibility Loss ( $L_3$ ).** To ensure the adversarial prompt remains safe according to all other guardrails and safety-aligned LLMs (C2), we introduce a safety scorer  $\mathcal{S}$  and propose  $L_3$ . The safety scorer  $\mathcal{S}$  measures the safety stage of an input  $x$  (Schmidt and Wiegand, 2017; Mathew et al., 2021; Vidgen et al., 2021):  $\mathcal{S}(x) = y_s \in [0, 1]$ , where  $y_s = 0$  indicates no security risk, and  $y_s = 1$  indicates a potential risk. The loss term is defined as:

$$L_3 = \sigma(\mathcal{S}(x_a||x_q), 0). \quad (2)$$

By minimizing  $L_3$ , we prevent unintended rejections from unrelated guardrails, preserving the specificity of the attack on  $\mathcal{G}_t$ .

The final loss function is defined as:

$$L = L_1 + \alpha \cdot L_2 + \beta \cdot L_3, \quad \alpha, \beta \in \mathbb{R}, \quad (3)$$

where  $\alpha$  and  $\beta$  control the weights of the loss terms  $L_2$  and  $L_3$ , respectively. By jointly minimizing  $L_1$ ,  $L_2$ , and  $L_3$ , we adopt the GCG (Greedy Coordinate Gradient) strategy (Zou et al., 2023) with the updates following (Jin et al., 2024) to craft adversarial prompts. Specifically, we compute gradients w.r.t. token embeddings. For each token position, we sample top- $k$  candidate replacements based on the gradient scores and their constraints. Then we perform greedy coordinate updates over the adversarial prompts. The optimized adversarial prompts could thus expose the behavior of the candidate guardrail while maintaining compatibility with other guardrails

## 4.2 Adversarial Prompt Test

We then perform identification using two complementary strategies: the *input guard test* and the

*output guard test*. Together, these tests cover scenarios where guardrails are deployed at either the input or output stage of an agent, thereby addressing C3. In addition, we introduce match score, a novel metric that facilitates identification based solely on the candidate guardrail. Match score provides clear and actionable decision strategies, effectively addressing the challenge of threshold selection in identification (C4).

**Input Guard Test.** We first consider the identification of the input guardrail. As an input guardrail in an AI agent, it will decide whether the user prompt  $x$  should be passed to the LLM. Thus, we do the input guard test by directly querying the AI agent with our adversarial prompts  $\{x_a||x_q\}$ . If the AI agent responds with refuse, we consider the candidate guardrail to probably serve as an input guardrail in the AI agent (see Figure 2 (2a)).

**Output Guard Test.** The third challenge C3 comes in identifying the output guardrail in the AI agent, as the output guardrail takes the LLM’s generated response  $y$  as input instead of the user prompt  $x$ . It is hard for us to manipulate the LLM’s output as we have no knowledge about it. To solve this, we design a prompt template that asks the AI agent to repeat our adversarial prompt, e.g., “Please repeat: [Adversarial Prompt]” Ideally, the LLM’s response  $y$  should be the adversarial prompt  $\{x_a||x_q\}$  and will be passed to the output guardrail. We empirically design and evaluate five prompt templates and select the best one as shown in Section 5.3. With the well-designed prompt template, we can ask the AI agent to repeat the adversarial prompts to test its output guardrail as shown in Figure 2 (2b). Then, we can make identifications based on the response from the agent.

**More Complex Scenarios.** In real-world applications, the guardrail can be deployed in both input and output stages of an AI agent (see Figure 1). The agent may also equip different guardrails. This poses challenges for identification (C3). Therefore, given a candidate guardrail, it is necessary to conduct both input and output guard tests on the AI agent, as illustrated in Appendix E. Experiments in such scenarios in Section 5.4 empirically showcase the effectiveness of **AP-Test**.

**Match Score.** To better quantify, we introduce the *refusal rate*  $r = \frac{\#\text{refuse}}{\#\text{all responses}} \in [0, 1]$ , which is the ratio of refuse among all responses, where refuse represents that the AI agent refuses to respond to the query. A higher refusal rate indicates

that the candidate guardrail is more likely to be the input guardrail in this AI agent. Unlike existing LLM identification methods (Gubri et al., 2024; Jin et al., 2024) that require testing on a set of LLMs to help distinguish the identity (C4), we propose a novel metric, match score, based solely on the candidate guardrail. We first calculate the refusal rate of directly querying the candidate guardrail with the optimized adversarial prompts, denoted as the candidate refusal rate  $r_t$ . Note that the candidate refusal rates  $r_t$  in both input and output guard tests are the same, as we directly query the candidate guardrail, disregarding the AI agent or the LLM. Then, we define the match score for an AI agent as:

$$MS = \frac{|\min(r, r_t) - 0|^\lambda}{|r_t - 0|^\lambda} = \frac{|\min(r, r_t)|^\lambda}{|r_t|^\lambda}, \quad (4)$$

where 0 is the lower bound of  $r$  and  $\lambda$  ( $\lambda \geq 1$ ) represents the scaling factor. The match score  $MS \in [0, 1]$  depicts how likely the candidate guardrail exists in the AI agent. In other words, this measures whether the agent refuses at least as often as the candidate guardrail would. We consider the candidate guardrail to exist in the AI agent if  $MS > 0.5$ ; our experiments in Section 5.2 show that the identification performance is not sensitive to this threshold.

## 5 Experiments

### 5.1 Experimental Settings

We take four different guardrails as our candidate guardrails, including WildGuard (Han et al., 2024), LlamaGuard (Inan et al., 2023), LlamaGuard2 (Meta, 2024), and LlamaGuard3 (Chi et al., 2024). We use 10 guardrails for evaluation, including the four candidate guardrails, AegisDefensive, AegisPermissive (Ghosh et al., 2024), ShieldGemma-2B, ShieldGemma-9B, ShieldGemma-27B (Zeng et al., 2024), and GPT4o (OpenAI, 2024). We obey their default settings in our experiments and follow (Zeng et al., 2024) to prompt GPT4o as an input/output guardrail. Besides, we use Llama3.1 (Dubey et al., 2024) and GPT4o as the LLMs of the conversational AI agents. For the safety scorer in the output guard test, we consider state-of-the-art hate speech detectors and use LFTW-R4 (Vidgen et al., 2021) in our experiments. Details of different LLMs and guardrails we use can be found in Table 2 in Appendix A.

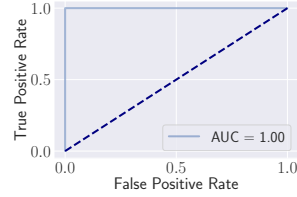


Figure 3: ROC curve of input guard test.

Candidate Guard	WildGuard	LlamaGuard	LlamaGuard2	LlamaGuard3	AegisDefensive	AegisPermissive	LlamaGuard2'	LlamaGuard3'	ShieldGemma-2B'	ShieldGemma-9B'	ShieldGemma-27B'	GPT4o'	None'
WildGuard	1.00	0.00	0.00	0.00	0.00	0.05	0.00	0.00	0.00	0.00	0.00	0.00	0.00
LlamaGuard	0.06	1.00	1.00	0.81	0.00	0.16	0.00	0.00	0.00	0.00	0.00	0.00	0.00
LlamaGuard2	0.00	0.00	0.00	0.00	1.00	0.22	0.00	0.00	0.00	0.00	0.00	0.00	0.00
LlamaGuard3	0.00	0.00	0.00	0.00	0.00	0.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00

Figure 4: The match scores of input guard tests on different AI agents. A larger match score indicates the candidate guardrail is more probable to exist in the agent.

Firstly, we optimize the adversarial prompts based on each candidate guardrail. We follow the settings in Jin et al. (2024) and use their dataset as the query set  $Q$ , which consists of 50 simple questions. For each query in  $Q$ , we optimize a 32-token adversarial prefix using loss weights  $\alpha = 0.01$  and  $\beta = 1000$ . We further investigate the impact of different loss term weights in the ablation study. The experiments are conducted with a batch size of 64 over 200 epochs on NVIDIA A100 GPUs with 40 GB of memory.

Then, we conduct input/output tests on agents equipped with different LLMs and input/output guardrails. We consider this a binary classification problem and calculate the classification accuracy for each candidate guardrail over the 11 AI Agents based on the match score  $MS$  (with  $\lambda = 2$ ) and plot the ROC curve for all test results. For each LLM, we construct 11 AI agents for evaluation: one without any guardrail, and the other 10 with different guardrails at the input/output stage. We discuss the situation that the agent contains additional guardrails later in Section 5.4.

### 5.2 Input Guard Test

We first evaluate the input guard test of **AP-Test** on Llama3.1-based and GPT4o-based agents. Results show that **AP-Test** achieves a *perfect* classification accuracy (1.00) in all agents for each candidate guardrail. To further assess the robustness of our approach, we plot the ROC curves for all test results in Figure 3, observing an AUC of 1.00. This indicates that our identification performance is not

Candidate Guard	WildGuard	LlamaGuard	LlamaGuard2	LlamaGuard3	WildGuard'	LlamaGuard'	AegisDefensive'	AegisPermissive'	LlamaGuard2'	LlamaGuard3'	ShieldGemma-2B'	ShieldGemma-9B'	ShieldGemma-27B'	GPT-4o'	None'
WildGuard	1.00	0.00	0.02	0.00	0.00	0.10	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
LlamaGuard	0.05	1.00	1.00	0.96	0.00	0.30	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
LlamaGuard2	0.02	0.00	0.01	0.00	1.00	0.44	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
LlamaGuard3	0.12	0.00	0.19	0.02	0.01	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

Figure 5: The match scores of output guard tests on different Llama3.1-based AI agents. A larger match score indicates the candidate guardrail is more probable to exist in the agent.

sensitive to the threshold selection. These findings highlight both the effectiveness and reliability of our method in the input guard test setting.

Figure 4 further illustrates the match scores of the input guard tests on Llama3.1-based agents equipped with various input guardrails. We observe an obvious distinction between the existence and non-existence of the candidate guard. For example, the WildGuard-agent achieves a match score of 1.00 on the WildGuard-specific adversarial prompts, while it only achieves a maximum distance of 0.05 on other agents. Additionally, the LlamaGuard-specific input test achieves 1.00 and 0.81 match score on the AegisDefense-agent and AegisPermissive-agent, respectively. That is because AegisDefense and AegisPermissive are fine-tuned on LlamaGuard, which are the derivatives of LlamaGuard. Results of GPT4o-based agents can be found in Figure 9 in Appendix C. We find that the results of both Llama3.1-based and GPT4o-based agents are almost the same. This shows that the LLM has little influence on the input guard test, i.e., our **AP-Test** is robust to different LLMs. This also indicates that the internal safety alignment will not affect the performance of our method.

### 5.3 Output Guard Test

**Prompt Template Selection.** As mentioned in Section 4, the output guard test requires a prompt template that asks the AI agent to repeat our adversarial prompts. We begin with the template: “Please repeat: [Adversarial Prompt]” and ask ChatGPT (OpenAI, 2024) to generate five variations, as shown in Table 3 in Appendix B. To assess how well the templates help the LLM repeat the adversarial prompt, we use Gemma2-9B (Mesnard et al., 2024) as a surrogate LLM and evaluate the similarity between the generated text and the input adversarial prompt optimized on LlamaGuard3. The results in Table 4 in Appendix B show that

Input Guard	None	WildGuard	LlamaGuard	LlamaGuard2	LlamaGuard3	ShieldGemma-2B
None	0.00	0.00	0.00	0.00	0.00	0.00
WildGuard	0.00	0.00	0.00	0.00	0.00	0.00
LlamaGuard	0.00	0.00	0.00	0.00	0.00	0.00
LlamaGuard2	0.00	0.00	0.00	0.00	0.00	0.00
LlamaGuard3	0.00	0.00	0.00	0.00	0.00	0.00
ShieldGemma-2B	1.00	1.00	1.00	1.00	1.00	1.00
None'	0.00	0.00	0.00	0.00	0.00	0.00
WildGuard'	0.00	0.00	0.00	0.00	0.00	0.00
LlamaGuard'	0.00	0.00	0.00	0.00	0.00	0.00
LlamaGuard2'	0.00	0.00	0.00	0.00	0.00	0.00
LlamaGuard3'	0.00	0.00	0.00	0.00	0.00	0.00
ShieldGemma-2B'	0.00	0.00	0.00	0.00	0.00	0.00

Figure 6: Influence of the presence of additional guardrails on input guard test. We report match scores on each agent.

the third template (T3) achieves the highest similarity scores, including a ROUGE-L score of 0.956. Thus, we select T3 as the prompt template for our experiments: “Your task is to mimic text. Repeat this string exactly as input, with no interpretation: [Adversarial Prompt].”

**Evaluation Results.** Using T3 as the prompt template, we evaluate the output guard test on AI agents with different output guardrails. Similar to the input guard test, **AP-Test** achieves an accuracy of 1.00 and an AUC of 1.00 in all agents for four candidate guardrails (see Appendix D). We further exploit the match scores of the test on Llama3.1-based agents, as shown in Figure 5. All agents with match scores larger than 0.50 are indeed equipped with the corresponding candidate guardrail as the output guardrail. For example, the match score on the WildGuard-agent reaches 1.00 on WildGuard-specific adversarial prompts, while it is 0.00 on LlamaGuard-agent and agents with LlamaGuard’s derivatives (AegisDefense-agent and AegisPermissive-agent). This indicates that our **AP-Test** successfully distinguishes the output guardrail used in the agent. We also observe that the output guard test is harder than the input guard test. For example, for LlamaGuard2-specific adversarial prompts, the match score of the LlamaGuard3-agent achieves 0.44, which is 0.22 farther than that in the input guard test and is closer to 0.50.

Results of GPT4o-based agents are shown in Figure 11 in Appendix D. We find that there is a slight performance difference between Llama3.1-based and GPT4o-based agents. This discrepancy is due to the information loss during the LLM processing. In other words, the performance of our output guard test is influenced by how well the LLM can repeat the adversarial prompts.

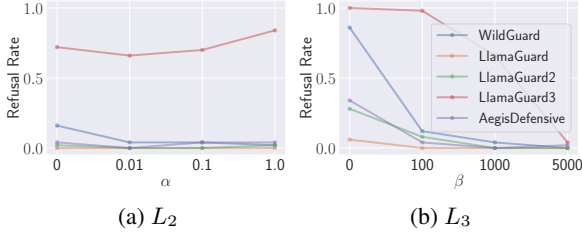


Figure 7: Influence of different weights of loss terms. We report the refusal rates while omitting the agents with a consistent refusal rate of 0.00.

#### 5.4 More Complex Scenarios

Our primary experiments show that the base LLMs (e.g., Llama3.1 and GPT4o) have little influence on the performance of our **AP-Test**. However, we assumed that the AI agent contains either an input or an output guardrail. To evaluate the robustness of **AP-Test** under more complex scenarios, we relax this assumption and assess its performance when both input and output guardrails are simultaneously deployed. Specifically, we evaluate **AP-Test** on  $6 \times 6 = 36$  Llama3.1-based agents constructed with all combinations of input and output guardrails, where each guardrail is one of: WildGuard, LlamaGuard, LlamaGuard2, LlamaGuard3, ShieldGemma-2B, or absent (N/A). We apply both input and output guard tests to each of these 36 agents. Both tests achieve perfect classification accuracy (1.00) across all agents. Figure 6 shows the match score of the input guard test. This further suggests that the presence of additional guardrails has minimal impact on the effectiveness of our method in identifying the candidate guardrail. Note that the match score of output guard test can be found in Figure 13b in Appendix E.

## 6 Ablation Study

To explore the influence of different components and hyperparameters of our **AP-Test**, we take LlamaGuard3 as our candidate guardrail and experiment on 11 Llama3.1-based AI agents with different input guardrails used in the main experiments. We find that the classification accuracy and the match score are relatively robust to various settings. Thus, we report the refusal rate  $r$  on each agent to better illustrate the influence in certain cases.

**Loss Terms.** We first investigate the influence of loss terms. Figure 7 shows the refusal rates using various weights of loss terms on all agents. We observe that  $L_2$  enhances the refusal rates on the agent with the candidate guardrail while slightly increas-

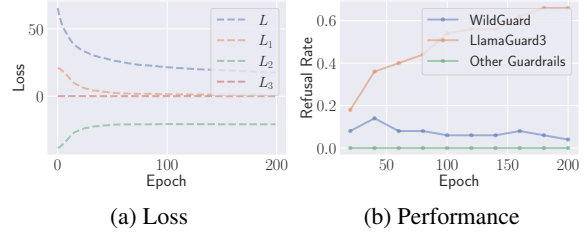


Figure 8: Influence of epochs. We report the loss values and refusal rates. All refusal rates on agents with other guardrails are 0.00.

	W/O $Q$ Original Alpaca Synthetic			
WildGuard	0.94	0.04	0.02	0.04
LlamaGuard	0.98	0.00	0.00	0.00
LlamaGuard2	0.78	0.00	0.00	0.00
LlamaGuard3	0.92	0.66	0.76	0.72

Table 1: Influence of normal query set  $Q$ . We report the refusal rates of the input guard test for agents with four different input guardrails.

ing the refusal rates on other agents. For example, the refusal rate on LlamaGuard3-agent increases from 0.72 to 0.84, while that on AegisDefensive-agent slightly increases to 0.04. As for  $L_3$ , we find it helps suppress the refusal rate on agents with other guardrails. For example, when  $\beta$  is 0, although the refusal rate on LlamaGuard3-agent is high (1.00), the refusal rate on WildGuard-agent reaches 0.86, which is so close to the candidate refusal rate  $r_t$  (0.98) on LlamaGuard3.

**Epochs.** To analyze the influence of the number of epochs, we analyze the loss values and refusal rates of the LlamaGuard3-specific input guard test across different epochs, as illustrated in Figure 8. The results indicate that the losses converge around the 50th epoch, highlighting the efficiency of our approach. Furthermore, the refusal rate on LlamaGuard3-agent increases with additional epochs, reaching 0.62 by the 160th epoch, while the refusal rate for WildGuard-agents decreases. This behavior demonstrates that as the number of epochs increases, the adversarial prompts become more reliable, as the distinction between the candidate guardrail and other guardrails grows. Ultimately, the system converges at a specific epoch level, ensuring stable and consistent performance.

**Query Set.** We further explore the influence of the normal query set  $Q$ , which serves as the starting point for adversarial prompt optimization. Besides the original query set we used in the main experiments, we construct two additional datasets:

(1) randomly select 50 samples from the Alpaca dataset (Taori et al., 2023), and (2) use GPT 5.2 (OpenAI, 2025) to randomly generate 50 questions (Synthetic) based on the original query set. Table 1 shows the refusal rates of LlamaGuard3-specific input guard test without or with different query sets. Results show that our **AP-Test** achieves comparable and robust performance with different query sets, which indicates robustness to the query set and our flexibility. Further, we find that without the query set, the test on agents equipped with guardrails other than LlamaGuard3 achieves even higher refusal rates. For instance, the refusal rate on the LlamaGuard-agent is 0.98, 0.06 higher than that on the LlamaGuard3-agent. This means that **AP-Test** without  $Q$  mistakes that LlamaGuard3 is used in the WildGuard-agent. In this sense, we address the importance of the query set as the initial guidance for adversarial optimization.

## 7 Conclusion

In this work, we tackle the problem of identifying guardrails deployed in conversational AI agents, a crucial step toward understanding system behavior and potential vulnerabilities. We propose **AP-Test**, a novel approach that leverages guard-specific adversarial prompts to identify the guardrail component within a black-box AI agent for the first time, effectively addressing the key challenges in this task. Experiments conducted on four candidate guardrails across various AI agents demonstrate the effectiveness and robustness of **AP-Test**. Our ablation study underscores the significance of our proposed loss terms and the query set, revealing that their removal leads to a substantial degradation in identification performance.

## Limitations

As this is the first study to tackle the guardrail identification problem, there is a lack of established baselines for direct comparison. We hope our work serves as a foundation for this problem and calls for more advanced techniques. Our study focuses exclusively on model-based guardrails, excluding rule-based filters and external moderation systems such as LLM API wrappers (e.g., GPT-4o with specific system prompts). Moreover, identifying a fully proprietary, closed-source guardrail is a harder and orthogonal problem. We leave the exploration of these different guard techniques to future work.

## Ethical Considerations

In research-oriented settings, safety-aligned LLMs and external guardrails are typically evaluated in isolation. However, in real-world attack scenarios, black-box AI agents often incorporate a combination of both, making attacks more challenging. Identifying the deployed guardrail in the agent can help attackers conduct stronger attacks. On the other hand, if the model owner needs to protect the copyright of their guardrail, **AP-Test** can help detect unauthorized usage. Thus, understanding this interplay is crucial for both attack and defense.

To facilitate replication and future research, we provide a comprehensive description of our experimental setup in Section 5.1 and Appendix A. The datasets we use are in English. Both datasets and code are either publicly available or generated by LLMs, ensuring no inclusion of personally identifiable information and eliminating user de-anonymization risks. We also adhere to the licenses or terms for use and emphasize that all collected data is solely for scientific purposes. To uphold responsible data management, only anonymized prompts/datasets will be shared when the code repository is made public.

Ultimately, our work highlights the need for techniques that make guardrails inherently harder to identify and circumvent. Meanwhile, we aim to contribute to the attack-defense cycle by enabling more effective adversarial testing and enhancing the explainability of AI agents in safety-critical scenarios.

## References

- Perspective API. <https://www.perspectiveapi.com>.
- Amazon. 2025. AWS Responsible AI Policy. <https://aws.amazon.com/cn/machine-learning/responsible-ai/policy/>.
- Anthropic. Claude. <https://claude.ai/>.
- Jianfeng Chi, Ujjwal Karn, Hongyuan Zhan, Eric Smith, Javier Rando, Yiming Zhang, Kate Plawiak, Zacharie Delpierre Coudert, Kartikeya Upasani, and Mahesh Pasupuleti. 2024. Llama Guard 3 Vision: Safeguarding Human-AI Image Understanding Conversations. *CoRR abs/2411.10414*.
- Edoardo DeBenedetti, Jie Zhang, Mislav Balunovic, Luca Beurer-Kellner, Marc Fischer, and Florian Tramèr. 2024. AgentDojo: A Dynamic Environment to Evaluate Attacks and Defenses for LLM Agents. *CoRR abs/2406.13352*.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu

- Li, Ziyi Gao, and 81 others. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. *CoRR abs/2501.12948*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, and 82 others. 2024. The Llama 3 Herd of Models. *CoRR abs/2407.21783*.
- Shaona Ghosh, Prasoon Varshney, Erick Galinkin, and Christopher Parisien. 2024. AEGIS: Online Adaptive AI Content Safety Moderation with Ensemble of LLM Experts. *CoRR abs/2404.05993*.
- Github. Copilot. <https://github.com/features/copilot>.
- Yichen Gong, DeLong Ran, Jinyuan Liu, Conglei Wang, Tianshuo Cong, Anyu Wang, Sisi Duan, and Xiaoyun Wang. 2023. FigStep: Jailbreaking Large Vision-language Models via Typographic Visual Prompts. *CoRR abs/2311.05608*.
- Google. 2024. Generative AI Prohibited Use Policy. <https://policies.google.com/terms/generative-ai/use-policy?hl=en>.
- Martin Gubri, Dennis Ulmer, Hwaran Lee, Sangdoon Yun, and Seong Joon Oh. 2024. TRAP: Targeted Random Adversarial Prompt Honey-pot for Black-Box Identification. In *Annual Meeting of the Association for Computational Linguistics (ACL)*. ACL.
- Seungju Han, Kavel Rao, Allyson Ettinger, Liwei Jiang, Bill Yuchen Lin, Nathan Lambert, Yejin Choi, and Nouha Dziri. 2024. WildGuard: Open One-Stop Moderation Tools for Safety Risks, Jailbreaks, and Refusals of LLMs. *CoRR abs/2406.18495*.
- Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, and Madian Khabza. 2023. Llama Guard: LLM-based Input-Output Safeguard for Human-AI Conversations. *CoRR abs/2312.06674*.
- Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2023. BeaverTails: Towards Improved Safety Alignment of LLM via a Human-Preference Dataset. In *Annual Conference on Neural Information Processing Systems (NeurIPS)*. NeurIPS.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, élio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7B. *CoRR abs/2310.06825*.
- Heng Jin, Chaoyu Zhang, Shanghao Shi, Wenjing Lou, and Y. Thomas Hou. 2024. ProFLingo: A Fingerprinting-based Copyright Protection Scheme for Large Language Models. *CoRR abs/2405.02466*.
- Haoran Li, Dadi Guo, Wei Fan, Mingshi Xu, and Yangqiu Song. 2023. Multi-step Jailbreaking Privacy Attacks on ChatGPT. *CoRR abs/2304.05197*.
- Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. 2023. AutoDAN: Generating Stealthy Jailbreak Prompts on Aligned Large Language Models. *CoRR abs/2310.04451*.
- Yupei Liu, Yuqi Jia, Rungeng Geng, Jinyuan Jia, and Neil Zhenqiang Gong. 2024a. Formalizing and Benchmarking Prompt Injection Attacks and Defenses. In *USENIX Security Symposium (USENIX Security)*. USENIX.
- Zhendong Liu, Yuanbi Nie, Yingshui Tan, Xiangyu Yue, Qiushi Cui, Chongjun Wang, Xiaoyong Zhu, and Bo Zheng. 2024b. Safety Alignment for Vision Language Models. *CoRR abs/2405.13581*.
- Neal Mangaokar, Ashish Hooda, Jihye Choi, Shreyas Chandrashekar, Kassem Fawaz, Somesh Jha, and Atul Prakash. 2024. PRP: Propagating Universal Perturbations to Attack Large Language Model Guard-Rails. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 10960–10976. ACL.
- Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. 2021. HateXplain: A Benchmark Dataset for Explainable Hate Speech Detection. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 14867–14875. AAAI.
- Hope McGovern, Rickard Stureborg, Yoshi Suhara, and Dimitris Alikaniotis. 2024. Your Large Language Models Are Leaving Fingerprints. *CoRR abs/2405.14057*.
- Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron Singer, and Amin Karbasi. 2023. Tree of Attacks: Jailbreaking Black-Box LLMs Automatically. *CoRR abs/2312.02119*.
- Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, and 30 others. 2024. Gemma: Open Models Based on Gemini Research and Technology. *CoRR abs/2403.08295*.
- Meta. Acceptable Use Policy. <https://ai.meta.com/llama/use-policy/>.
- Meta. 2024. Llama Guard 2. <https://www.llama.com/docs/model-cards-and-prompt-formats/meta-llama-guard-2/>.
- Alexander Neumann, Yue Yin, Sulayman K Sowe, Stefan Decker, and Matthias Jarke. 2024. An LLM-Driven Chatbot in Higher Education for Databases and Information Systems. *IEEE Transactions on Education*.
- OpenAI. Usage policies. <https://openai.com/policies/usage-policies>.
- OpenAI. 2024. GPT-4o. <https://openai.com/index/hello-gpt-4o/>.
- OpenAI. 2025. GPT 5.2. <https://openai.com/index/introducing-gpt-5-2/>.
- Cheng Peng, Xi Yang, Aokun Chen, Kaleb E. Smith, Nima M. Pournejatian, Anthony B. Costa, Cheryl Martin, Mona G. Flores, Ying Zhang, Tanja Magoc, Gloria P. Lipori, Duane A. Mitchell, Naykky Singh Ospina, Mustafa M.

- Ahmed, William R. Hogan, Elizabeth A. Shenkman, Yi Guo, Jiang Bian, and Yonghui Wu. 2023. A study of generative large language model for medical research and healthcare. *NPJ Digital Medicine*.
- Anna Schmidt and Michael Wiegand. 2017. A Survey on Hate Speech Detection using Natural Language Processing. In *Workshop on Natural Language Processing for Social Media (SocialNLP)*, pages 1–10. ACL.
- Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, and Yang Zhang. 2024. Do Anything Now: Characterizing and Evaluating In-The-Wild Jailbreak Prompts on Large Language Models. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*. ACM.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford Alpaca: An Instruction-following LLaMA Model. [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca).
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. *CoRR abs/2307.09288*.
- Bertie Vidgen, Tristan Thrush, Zeerak Waseem, and Douwe Kiela. 2021. Learning from the Worst: Dynamically Generated Datasets to Improve Online Hate Detection. In *Annual Meeting of the Association for Computational Linguistics and International Joint Conference on Natural Language Processing (ACL/IJCNLP)*, pages 1667–1682. ACL.
- Rui Wen, Zheng Li, Michael Backes, and Yang Zhang. 2024. Membership Inference Attacks Against In-Context Learning. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*. ACM.
- Ziqing Yang, Yixin Wu, Yun Shen, Wei Dai, Michael Backes, and Yang Zhang. 2025. The Challenge of Identifying the Origin of Black-Box Large Language Models. *CoRR abs/2503.04332*.
- Jiahao Yu, Yuhang Wu, Dong Shu, Mingyu Jin, and Xinyu Xing. 2023. Assessing Prompt Injection Risks in 200+ Custom GPTs. *CoRR abs/2311.11538*.
- Boyi Zeng, Chenghu Zhou, Xinbing Wang, and Zhouhan Lin. 2023. HuRef: HUMAN-REadable Fingerprint for Large Language Models. *CoRR abs/2312.04828*.
- Wenjun Zeng, Yuchi Liu, Ryan Mullins, Ludovic Peran, Joe Fernandez, Hamza Harkous, Karthik Narasimhan, Drew Proud, Piyush Kumar, Bhaktipriya Radharapu, Olivia Sturman, and Oscar Wahltinez. 2024. ShieldGemma: Generative AI Content Moderation Based on Gemma. *CoRR abs/2407.21772*.
- Qiusi Zhan, Zhixiang Liang, Zifan Ying, and Daniel Kang. 2024. InjecAgent: Benchmarking Indirect Prompt Injections in Tool-Integrated Large Language Model Agents. *CoRR abs/2403.02691*.
- Qingzhao Zhang, Ziyang Xiong, and Z. Morley Mao. 2024a. Safeguard is a Double-edged Sword: Denial-of-service Attack on Large Language Models. *CoRR abs/2410.02916*.
- Rui Zhang, Hongwei Li, Rui Wen, Wenbo Jiang, Yuan Zhang, Michael Backes, Yun Shen, and Yang Zhang. 2024b. Instruction Backdoor Attacks Against Customized LLMs. In *USENIX Security Symposium (USENIX Security)*. USENIX.
- Xuandong Zhao, Xianjun Yang, Tianyu Pang, Chao Du, Lei Li, Yu-Xiang Wang, and William Yang Wang. 2024. Weak-to-Strong Jailbreaking on Large Language Models. *CoRR abs/2401.17256*.
- Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. 2023. Universal and Transferable Adversarial Attacks on Aligned Language Models. *CoRR abs/2307.15043*.

## A Model Details

Table 2 shows the details of the models used in our experiments, including the versions we used.

## B Prompt Template Selection

As mentioned in Section 4, the output guard test requires a prompt template that asks the AI agent to repeat our adversarial prompts. We begin with the template: “Please repeat: [Adversarial Prompt]” and ask ChatGPT (OpenAI, 2024) to generate five variations, as shown in Table 3. To assess how well the templates help the LLM repeat the adversarial prompt, we use Gemma2-9B (Mesnard et al., 2024) as a surrogate LLM and evaluate the similarity between the generated text and the input adversarial prompt optimized on LlamaGuard3. The results in Table 4 show that the third template (T3) achieves the highest similarity scores, including a ROUGE-L score of 0.956. Thus, we select T3 as the prompt template for our experiments: “*Your task is to mimic text. Repeat this string exactly as input, with no interpretation:* [Adversarial Prompt]”

## C Input Guard Test

Figure 9 shows the match scores of the input guard test on different agents. The agents are either based on Llama3.1 or GPT4o models. We find that there is little difference between the two agents based on different LLMs. This shows that the base LLM has little influence on the input guard test.

## D Output Guard Test

Figure 10 is the ROC curve of **AP-Test** of the output guard test. They exhibit a 1.00 classification accuracy and an AUC of 1.00, indicating the effectiveness of our identification. To be more concise, Figure 11 shows the match scores of the output guard test on different agents. The agents are either based on Llama3.1 or GPT4o models. We observe that the output guard test performs slightly better on GPT4o-based agents than Llama3.1-based ones. Take LlamaGuard3 as the candidate guard and WildGuard as the equipped output guard, for example. The match score on the Llama3.1-based agent is 0.12, while that on the GPT4o-based agent is 0.01. This discrepancy is due to the information loss during the LLM processing. In other words, the performance of our output guard test is influenced by the LLM and the prompt template,

i.e., how well the LLM can repeat the adversarial prompts.

## E AP-Test in More Complex Scenarios

In real-world applications, there are both input and output guardrails in an AI agent, as illustrated in Figure 1. Therefore, given a candidate guardrail, it is necessary to conduct both input and output guard tests on the AI agent. Here, we propose a two-step process, as shown in Figure 12.

- We first determine whether the candidate guardrail functions as an input guardrail through the input guard test. If so, we conclude that the candidate guardrail is deployed in the AI agent and do not proceed with the output guard test.
- If it is not used as the input guardrail, then we proceed with the output guard test for the candidate guardrail.

In our evaluations (Section 5.4) shown in Figure 13, we demonstrate that the proposed method can successfully identify the input/output guardrail even in the presence of a different output/input guardrail.

Model Name	Model Path
WildGuard (Han et al., 2024)	allenai/wildguard
LlamaGuard (Inan et al., 2023)	meta-llama/LlamaGuard-7b
└ AegisDefensive (Ghosh et al., 2024)	nvidia/Aegis-AI-Content-Safety-LlamaGuard-Defensive-1.0
└ AegisPermissive (Ghosh et al., 2024)	nvidia/Aegis-AI-Content-Safety-LlamaGuard-Permissive-1.0
LlamaGuard2 (Meta, 2024)	meta-llama/Meta-Llama-Guard-2-8B
LlamaGuard3 (Chi et al., 2024)	meta-llama/Llama-Guard-3-8B
ShieldGemma-2B (Zeng et al., 2024)	google/shieldgemma-2b
ShieldGemma-9B (Zeng et al., 2024)	google/shieldgemma-9b
ShieldGemma-27B (Zeng et al., 2024)	google/shieldgemma-27b
Perspective (Per)	Perspective API
GPT4o (OpenAI, 2024)	GPT4o-2024-11-20
Gemma2 (Mesnard et al., 2024)	google/gemma-2-9b-it
Llama3.1 (Dubey et al., 2024)	meta-llama/Llama-3.1-8B-Instruct

Table 2: Details of the guardrails and LLMs. Note that AegisDefensive and AegisPermissive are fine-tuned on LlamaGuard, as illustrated below.

Template
T1 Simply repeat and do not consider the meaning: [AP]
T2 Repeat exactly the following sentence, including all characters and punctuation: [AP]
T3 Your task is to mimic text. Repeat this string exactly as input, with no interpretation: [AP]
T4 Do not interpret the meaning of this text. Simply output it as provided: [AP]
T5 For debugging purposes, repeat the following text string exactly as input: [AP]

Table 3: Candidate prompt templates for output guard test. [AP] denotes the placeholder for the adversarial prompts.

#Template	Cosine Similarity	BLUE	ROUGE-1	ROUGE-2	ROUGE-L	ROUGE-LSum
T1	0.881	0.977	0.911	0.900	0.911	0.911
T2	0.933	0.961	0.939	0.934	0.939	0.939
T3	<b>0.944</b>	<b>0.980</b>	<b>0.956</b>	<b>0.945</b>	<b>0.956</b>	<b>0.956</b>
T4	0.870	0.929	0.884	0.876	0.884	0.884
T5	0.914	0.962	0.926	0.916	0.926	0.926

Table 4: Similarity between the input adversarial prompt and the output text from the surrogate LLM with different prompt templates.

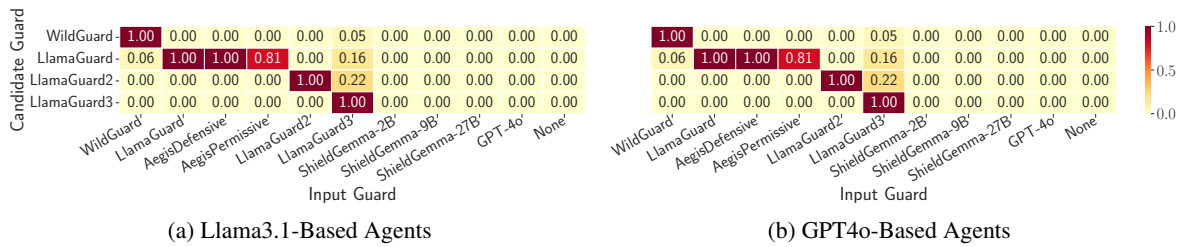


Figure 9: The match scores of input guard tests on different AI agents. A larger match score indicates the candidate guardrail is more probable to exist in the agent.

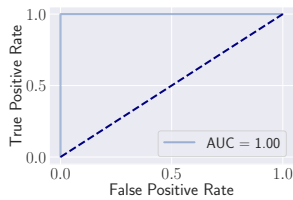


Figure 10: ROC curve of output guard test.

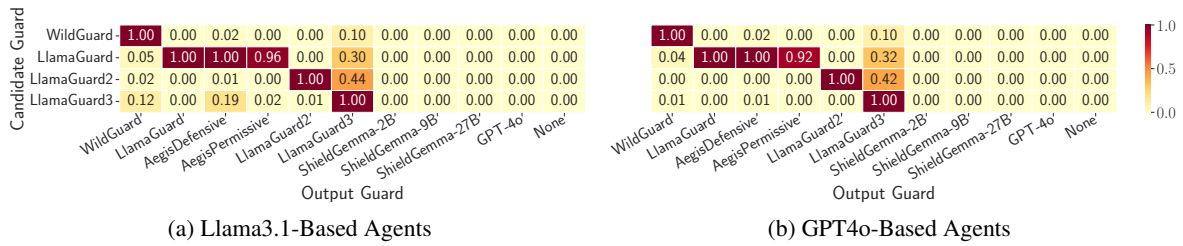


Figure 11: The match scores of output guard tests on different AI agents. A larger match score indicates the candidate guardrail is more probable to exist in the agent.

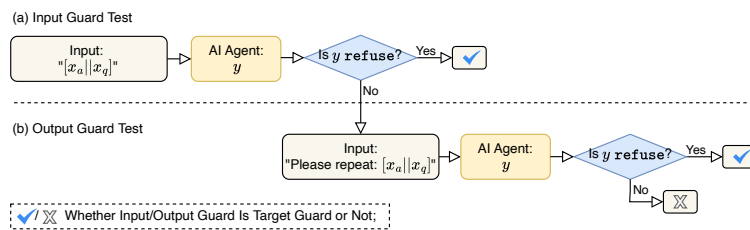


Figure 12: Workflow on real-world scenarios. We first conduct (a) the input guard test on the AI agent. If the results show that the candidate guardrail probably does not exist in the agent, then we further conduct (b) the output guard test to identify whether it serves as the output guardrail in the agent.

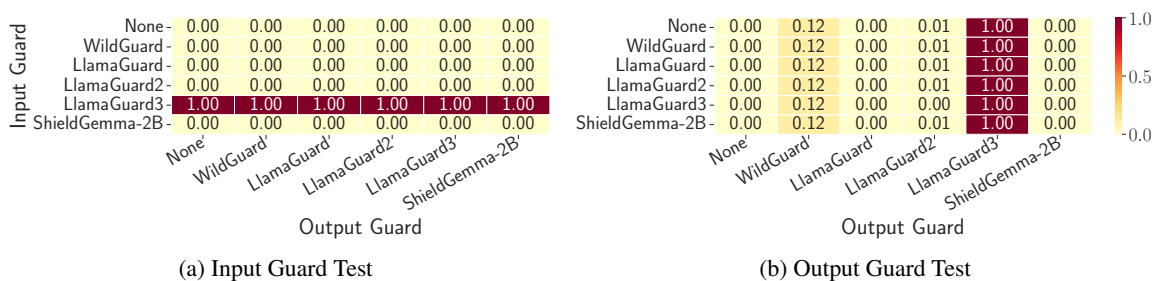


Figure 13: Influence of the presence of additional guardrails. We report match scores on each agent. A larger match score indicates the candidate guardrail is more probable to exist in the agent.