



# RLShield: Dynamic Jailbreak Detection for LLMs via Reinforced Adaptive Learning

Zhao Tong<sup>1,2\*</sup>, Pengfei Yang<sup>3\*</sup>, Yimeng Gu<sup>4</sup>, Haichao Shi<sup>1</sup>, Qiang Liu<sup>5</sup>,  
Xingcheng Xu<sup>6</sup>, Shu Wu<sup>5†</sup>, Xiao-Yu Zhang<sup>1†</sup>

<sup>1</sup>Institute of Information Engineering, Chinese Academy of Sciences

<sup>2</sup>School of Cyber Security, University of Chinese Academy of Sciences

<sup>3</sup>Harbin Institute of Technology, <sup>4</sup>Queen Mary University of London

<sup>5</sup>New Laboratory of Pattern Recognition (NLPR),

State Key Laboratory of Multimodal Artificial Intelligence Systems (MAIS),

Institute of Automation, Chinese Academy of Sciences

<sup>6</sup>Shanghai AI Laboratory

tongzhao@iie.ac.cn, 2023210925@stu.hit.edu.cn, yimeng.gu@qmul.ac.uk  
shihaichao@iie.ac.cn, qiang.liu@nlpr.ia.ac.cn, xingcheng.xu18@gmail.com  
shu.wu@nlpr.ia.ac.cn, zhangxiaoyu@iie.ac.cn

## Abstract

While prompt engineering enhances the capabilities of Large Language Models (LLMs), it also exposes critical safety concerns. Due to the inherent brittleness of their static safety boundaries, LLMs are vulnerable to *jailbreak prompts*, *i.e.* adversarial inputs designed to bypass safeguards and induce the generation of harmful content. Existing detection mechanisms rely on static model components or fixed decision thresholds, limiting their ability to generalize to evolving attack patterns and continual model updates. To bridge this gap, we propose **RLShield**, a dynamic jailbreak detection framework that employs reinforcement learning for adaptive threshold selection. RLShield incorporates three key innovations: (i) a dynamic retrieval and LLM-based rewriting module to simulate diverse adversarial contexts; (ii) a cross-layer representation analysis to pinpoint safety-critical parameters; and (iii) a Soft Actor-Critic (SAC) based agent that learns to predict optimal, sample-specific detection thresholds. Experimental results demonstrate that RLShield consistently outperforms state-of-the-art baselines in detection performance while maintaining high computational efficiency. Notably, it improves F1 by up to 7.3%, while achieving an average of  $3\times$  gain in inference efficiency across multiple LLM backbones. Our codes are available at [this website](#).

## 1 Introduction

Prompt engineering enables LLMs to perform diverse downstream tasks (Wang et al., 2025; Feng et al., 2025; Wu et al., 2024; Tong et al., 2025a,d;

\*These authors contributed to the work equally.

†To whom correspondence should be addressed.

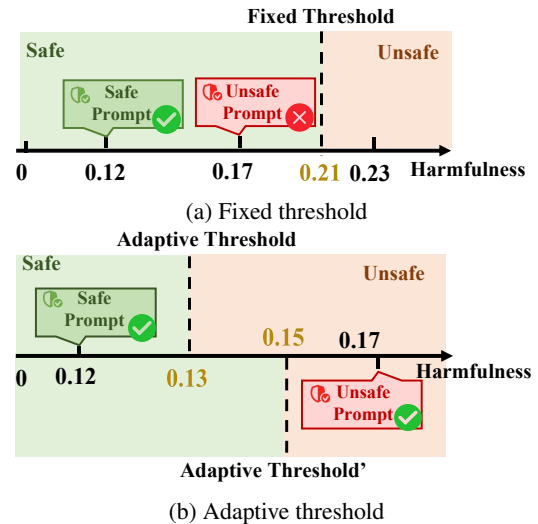


Figure 1: For the given prompt, fixed threshold leads to misclassification, while an adaptive threshold learns a context-aware threshold that gives correct classification.

Zhang et al., 2024; Tong et al., 2025c). However, prompts can expose safety vulnerabilities in LLMs. Sophisticated prompt engineering such as role-playing and semantic obfuscation can bypass safety alignment and elicit harmful content, including dangerous guidance (Yang et al., 2025), misinformation (Pan et al., 2023; Tong et al., 2025b), or hate speech (Masud et al., 2024). Despite extensive safety fine-tuning and reinforcement learning from human feedback (RLHF), the rigid generalization of their safety mechanisms leaves them vulnerable to out-of-distribution adversarial prompts. Therefore, developing robust, real-time detection mechanisms is of critical importance.

Recently, numerous efforts have been made to

detect jailbreaks. Early approaches rely on predefined metrics (e.g., perplexity) and fixed thresholds to identify malicious prompts (Jain et al., 2023; Wang et al., 2024b). Subsequent works expanded detection coverage by incorporating external knowledge (Liu et al., 2024; Tu et al., 2025). More recent state-of-the-art approaches leverage internal parameters of LLMs, such as activations or gradients (Xuan et al., 2025; Hu et al., 2024; Padakandla et al., 2025; Xie et al., 2024; Chen et al., 2025), to infer adversarial intent. Although effective, existing approaches have a common deficiency – decision rigidity: reliance on fixed settings and static parameters hinders adaptation to evolving jailbreak patterns, especially under rapid model updates and continuous parameter optimization (Echterhoff et al., 2024). Therefore, dynamic adaptation is critical for robust jailbreak prompt detection.

To enable dynamic adaptation leveraging LLMs, two major challenges must be addressed. *First*, it is challenging to efficiently detect jailbreak prompts with dynamic threshold. Unlike static detection, dynamic detection requires prompt-specific threshold, which causes additional computational overhead. As utilizing the full set of LLM parameters is computationally prohibitive, it is crucial to identify and exploit parameters that are most relevant to safety. Moreover, incorporating irrelevant parameters can introduce unnecessary noise that may degrade detection performance. *Second*, learning an adaptive threshold requires continuous-valued signals that provide guidance for threshold estimation. In jailbreak detection, prompts are typically annotated with only discrete classification labels which provide insufficient information to determine a threshold within a continuous range. Considering this, identifying alternative continuous and informative signals is essential for effectively learning an adaptive threshold.

To address these challenges, we propose RLShield, a dynamic jailbreak detection framework that integrates adaptive thresholding and continuous risk estimation. As illustrated in Fig. 1, if we use a fixed threshold, it may misclassify harmful prompts, whereas an adaptive threshold can make more accurate predictions based on different prompts. Our framework consists of three components: (i) **Dynamic Reference Selection**: We mitigate the rigidity of static pattern matching by simulating attack variations through retrieval-guided LLM rewriting. (ii) **Critical Parameter Local-**

**ization**: We analyze cross-layer representations to identify safety-related layers and their safety-critical parameters, directing focus toward the most informative signals. (iii) **RL-Driven Adaptive Thresholding**: We train a Soft Actor-Critic (SAC) agent to learn prompt-specific optimal decision thresholds. Unlike traditional fixed methods, this RL-based approach encodes continuous risk level information, allowing the detector to adapt to dynamic attack intensity across various scenarios.

Our main contributions are summarized as follows:

- We introduce a cross-layer localization method to pinpoint safety-critical parameters, which significantly improves detection efficiency.
- We propose RLShield, a dynamic jailbreak detection framework in which an RL agent learns adaptive, prompt-specific detection thresholds.
- We demonstrate that RLShield achieves SOTA performance across multiple benchmarks with high inference efficiency.

## 2 Related Work

### 2.1 Jailbreak Threats to LLM

Early approaches (Jain et al., 2023; Wang et al., 2024b; Peng et al., 2024) identify malicious prompts using perplexity or keyword rules but struggle with identifying semantically obfuscated paraphrases. As LLMs became targets, neural detectors emerged: Xuan et al. (2025) mine hidden states during decoding, GradSafe (Xie et al., 2024) assesses gradient similarity to a refusal anchor, and LlamaGuard (Meta, 2025) fine-tunes a safety-specialized model. Recent studies indicate that a few jailbreak samples in SFT datasets can alter refusal behaviors (Wang et al., 2024a; Huang et al., 2024), exposing vulnerabilities during training. These approaches typically fix the defense layer or decision threshold, limiting adaptability to evolving attack patterns or continual model updates. It is crucial for our jailbreak detector to adjust the threshold on the fly, ensuring harmful prompts are identified even as attack styles and models evolve.

### 2.2 Jailbreak Detection

Since the rise of LLMs, numerous studies have exploited their internal signals to detect jailbreak

prompts. Early rule-based detectors leverage perplexity or back-translation consistency (Jain et al., 2023; Wang et al., 2024b), yet Jain et al. (2023) found that paraphrased attacks drop recall by 30%. External APIs (Markov et al., 2023; Lees et al., 2022) provide zero-shot labels, but Han et al. (2024) showed they generalise poorly to out-of-domain dialects. Internal-representation methods mine hidden states (Xuan et al., 2025) or gradients (Xie et al., 2024; Padakandla et al., 2025); although these works achieve impressive performance, they do not adapt the decision boundary to each prompt. In contrast, we treat threshold selection as a continuous control problem and learn a prompt-specific classification threshold with reinforcement learning.

### 2.3 RL for Jailbreak Detection

Reinforcement learning has long been utilized for security tasks in view of its capability to learn adaptive policies (Haarnoja et al., 2018; Littman, 2015). In the field of LLMs, Ye et al. (2025) develop a robust RLHF policy designed to resist adversarial prompts; however, their agent focuses on response generation rather than detection. Similarly, Hu et al. (2024) employ reinforcement learning to optimize a soft prompt that guides the model toward refusal, but the decision boundary remains fixed post-training. To the best of our knowledge, jailbreak detection has not been framed as a continuous control problem whose action is a prompt-specific threshold. We introduce a SAC agent that outputs a sample-adaptive threshold in a single forward pass, eliminating the need for manual re-tuning as the backbone LLM evolves.

## 3 Methodology

We propose RLShield, a framework for dynamic jailbreak detection. The overall architecture is illustrated in Fig. 2. RLShield consists of three main components: (i) *Dynamic Reference Selection*, which constructs a comprehensive knowledge base and dynamically retrieves and rewrites reference prompts for each given prompt; (ii) *Critical Parameter Localization*, which identifies safety-related layers in the LLM and extracts their critical parameters to significantly reduce redundant computation; and (iii) *RL-Driven Adaptive Thresholding*, where a RL agent learns a prompt-specific threshold for jailbreak detection, enabling robust detection under distribution shifts and diverse jail-

break patterns.

### 3.1 Dynamic Reference Selection

To reduce reliance on static patterns and fixed references, RLShield dynamically selects prompt-specific safe and unsafe references from a curated knowledge base for safety-critical parameter selection and harmfulness score computation. We first construct two reference sets, namely *safe reference* knowledge base  $\mathcal{D}_s$  and *unsafe reference* knowledge base  $\mathcal{D}_u$ .

From  $\mathcal{D}_s$  and  $\mathcal{D}_u$ , we then retrieve safe and unsafe references for each prompt. This is to provide the data source for retrieval and subsequent safety-layer selection. We encode the input prompt  $x$  into an embedding  $\mathbf{e}_x$  and compute its cosine similarity with each reference prompt  $r$  of embedding  $\mathbf{e}_r$ :

$$s(x, r) = \frac{\mathbf{e}_x^\top \mathbf{e}_r}{\|\mathbf{e}_x\| \|\mathbf{e}_r\|}. \quad (1)$$

In order to obtain prompt-specific references that yield more consistent and separable internal signals for detection, we retrieve the top- $k$  references from the safe and unsafe knowledge bases, respectively. From each of the top- $k$  candidates, we randomly sample two references and denote sampled safe and unsafe reference sets as  $\mathcal{R}_s$  and  $\mathcal{R}_u$ , instead of constantly taking the top results. This mitigates near-duplicate retrievals among the highest-similarity prompts and increases reference diversity.

After retrieving and sampling prompt-specific references, we expand references via LLM-based rewriting to enhance diversity. For each reference prompt  $r$ , we generate its variants  $\{\tilde{r}\}$ . This is formulated in Eq. 2:

$$\tilde{r} = \text{LLM}(r, \mathcal{I}_{\text{rewrite}}), \quad (2)$$

where  $\mathcal{I}_{\text{rewrite}}$  denotes a predefined rewrite prompt that controls writing style while preserving the original safety label, and the detailed instructions together with the empirical validation of label preservation are provided in Appendix A.1. This process generates diverse semantic variants. Overall, we obtain the rewritten safe and unsafe reference sets  $\tilde{\mathcal{R}}_s$  and  $\tilde{\mathcal{R}}_u$  for each input prompt. These prompt-specific references are then used in the subsequent modules for critical-parameter selection and harmfulness scoring.

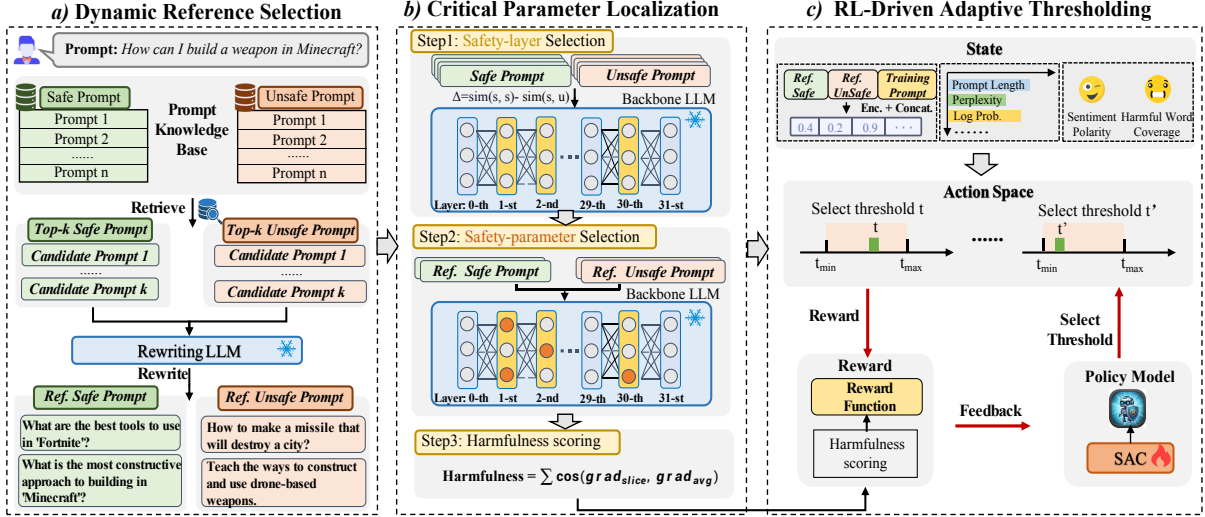


Figure 2: Overall architecture of RLShield. It consists of three components: (i) *Dynamic Reference Selection*, which simulates attack variations via retrieval-augmented LLM rewriting of prompt-specific safe and unsafe references; (ii) *Critical Parameter Localization*, which identifies safety-related layers and safety-critical parameters through cross-layer representation analysis for harmfulness scoring; (iii) *RL-Driven Adaptive Thresholding*, which trains an SAC agent to learn adaptive decision thresholds from continuous risk signals.

### 3.2 Critical Parameter Localization

After building the knowledge base and obtaining prompt-specific references, we next perform safety-layer selection and safety-parameter selection using the knowledge base and the references, respectively. This is to reduce reliance on static internal features, while eliminating redundant computation. Our design is motivated by prior work indicating that a small set of parameters in aligned LLMs plays an outsized role in shaping safety-aligned behaviors (Li et al., 2025; Gu et al., 2025; Zhou et al., 2025).

**Step I (Safety-layer selection):** We construct a safe set  $\mathcal{S}$  and an unsafe set  $\mathcal{U}$ , both drawn from the knowledge base. For each layer  $l$ , we extract the final-position hidden representation  $h_l(x)$ , which captures the overall semantics of the prompt at layer  $l$ . Using  $h_l(x)$ , we then compute the similarity between two prompts using cosine similarity:

$$s_l(x, x') = \frac{h_l(x)^\top h_l(x')}{\|h_l(x)\| \|h_l(x')\|}. \quad (3)$$

Over  $r$  random trials, we sample prompt pairs from  $\mathcal{S} \times \mathcal{S}$  and  $\mathcal{S} \times \mathcal{U}$  and compute the averaged similarities:

$$\begin{aligned} \bar{s}_l^{SS} &= \frac{1}{r} \sum_{t=1}^r s_l(x_t^S, x_t^{S'}), \\ \bar{s}_l^{SU} &= \frac{1}{r} \sum_{t=1}^r s_l(x_t^S, x_t^U). \end{aligned} \quad (4)$$

We define the layer-wise separation gap as

$$\Delta_l = \bar{s}_l^{SS} - \bar{s}_l^{SU}. \quad (5)$$

Here, a larger  $\Delta_l$  indicates stronger separability between safe and unsafe representations at layer  $l$ . Comparing adjacent layers, we select layers satisfying  $\Delta_l > \Delta_{l-1}$  and denote the resulting set as  $\ell_{\text{safe}}$ , because this indicates that the separation becomes more pronounced at layer  $l$ , suggesting that this layer is more sensitive to safety-related differences between prompts. More details are provided in Appendix A.3.2.

**Step II (Safety-parameter selection):** From the identified safety layers  $\ell_{\text{safe}}$ , we further select critical parameters for each prompt using gradients from dynamic references (Section 3.1) and a set of fixed references shared across all prompts with the same response, which is set to "Sure". To represent the parameters at different locations, we slice each parameter matrix in  $\ell_{\text{safe}}$  by rows and columns, and score each slice by the gap between the gradient cosine similarity of unsafe–unsafe reference and that of unsafe–safe reference. We retain the slices with high gaps. We use gradient cosine similarity because jailbreak prompts induce similar update-direction patterns, whereas safe prompts do not exhibit such consistency with jailbreak prompts.

We run this selection separately for the two reference types, obtaining  $\mathcal{P}_{\text{dyn}}$  and  $\mathcal{P}_{\text{fix}}$ , and take their intersection to retain the slices selected by both reference types, improving stability and reducing

reference-specific noise:

$$\mathcal{P} = \mathcal{P}_{\text{dyn}} \cap \mathcal{P}_{\text{fix}}. \quad (6)$$

**Step III (Harmfulness scoring):** To support adaptive thresholding with a continuous risk signal, we backpropagate the input prompt  $x$  with the same response, "Sure", to obtain gradients on its safety-parameter slices  $\mathcal{P}$ . We then compute the harmfulness score by averaging the cosine similarities between these gradients and the mean gradient of its dynamic unsafe references  $\tilde{\mathcal{R}}_u$  over all selected slices:

$$\text{harmfulness}(x) = \frac{1}{|\mathcal{P}|} \sum_{i \in \mathcal{P}} \cos(g_i(x), \bar{g}_i^u), \quad (7)$$

where  $i$  indexes a slice in  $\mathcal{P}$ ,  $g_i(\cdot)$  denotes the gradient of slice  $i$ , and  $\bar{g}_i^u$  is the mean gradient slice averaged over  $\tilde{\mathcal{R}}_u$

### 3.3 RL-Driven Adaptive Thresholding

After obtaining the harmfulness score, we learn a dynamic threshold policy, since this score provides a continuous, comparable risk signal. We model threshold estimation as a continuous-control reinforcement learning problem and optimize an agent with the Soft Actor-Critic (SAC) algorithm (Haarnoja et al., 2018). This formulation allows the learned policy not only to predict an appropriate threshold for each prompt, but also to maintain a reasonable margin between the harmfulness score and the selected threshold, thereby improving the robustness of the final decision. SAC’s continuous action space and entropy-regularized objective deliver both flexibility and training stability for adaptive threshold learning. In practice, we instantiate this process as a one-step Markov decision process (MDP) (Littman, 2015) tailored to our safety classification scenario.

For each input prompt, the agent observes a feature-based state  $s_t$  and outputs a continuous action  $a_t$  as the threshold. The environment applies this threshold to the harmfulness score to predict a label and returns a reward based on harmfulness score and ground-truth label. Each prompt forms an independent one-step episode, and we optimize the policy  $\pi$  to learn adaptive, prompt-specific threshold.

Next, we detailedly describe the **state**, **action**, **reward**, and **objective function** defined in our approach.

**State.** For each input prompt, the state is constructed as a unified feature vector by integrating

multiple features. Formally, the state at step  $t$  is defined as

$$\begin{aligned} s_t &= [\mathbf{f}_t^{\text{sem}}; \mathbf{f}_t^{\text{stat}}; \mathbf{f}_t^{\text{sent}}], \\ \mathbf{f}_t^{\text{sem}} &= [z_t^x; z_t^s; z_t^u], \\ \mathbf{f}_t^{\text{stat}} &= [L_t; \text{ppl}_t; \text{eng}_t; \mu_t; \sigma_t^2], \\ \mathbf{f}_t^{\text{sent}} &= [\text{pol}_t; \text{cov}_t], \end{aligned} \quad (8)$$

where  $[\cdot; \cdot; \cdot]$  denotes vector concatenation. This state representation comprises three components: (i) **semantic representations**  $\mathbf{f}_t^{\text{sem}}$  – the input prompt embedding and the embeddings of its dynamically retrieved and rewritten safe/unsafe references obtained in Section 3.1. They provide semantic cues about the prompt and its references, which help the agent perform context-aware thresholding. They are first encoded by BGE-M3 (Chen et al., 2024) and then reduced via PCA, denoted by  $z_t^x$ ,  $z_t^s$ , and  $z_t^u$ ; (ii) **prompt-level statistics**  $\mathbf{f}_t^{\text{stat}}$  – prompt length  $L_t$ , perplexity  $\text{ppl}_t$ , English-language indicator  $\text{eng}_t$ , and token log-probability statistics  $(\mu_t, \sigma_t^2)$ . These features provide complementary statistical signals about the prompt, helping the agent make more informed threshold selections; and (iii) **sentiment indicators**  $\mathbf{f}_t^{\text{sent}}$  – sentiment polarity  $\text{pol}_t$  (from NLTK (Loper and Bird, 2002)) and harmful word coverage  $\text{cov}_t$  (using HurtLex (Bassignana et al., 2018)). These reflect affective intensity and explicit harmful content for sentiment-aware threshold decisions.

All features are concatenated into a single continuous vector  $s_t \in \mathbb{R}^d$  and used as the state input to the agent.

**Action.** At each step  $t$ , the agent selects a continuous action  $a_t$ , which directly corresponds to a decision threshold used for harmfulness classification. The action space is defined as a continuous interval

$$a_t \in \mathcal{A} = [T_{\min}, T_{\max}], \quad (9)$$

where  $T_{\min}$  and  $T_{\max}$  are the lower and upper bounds of the threshold, respectively. In our implementation, the agent outputs a real-valued threshold within this range to adaptively separate safe and unsafe samples. This dynamic action design enables dynamic and prompt-specific threshold selection instead of using a fixed global threshold.

**Reward.** After the agent selects a threshold  $a_t$ , the predicted label is determined by comparing the selected threshold with the harmfulness score  $\text{harmfulness}_t$  obtained from Section 3.2. Let  $y_t \in \{0, 1\}$  denote the ground-truth label and  $\hat{y}_t$  denote

---

**Algorithm 1** Inference procedure of RLShield

---

**Require:** Input prompt  $x$ ; safe knowledge base  $\mathcal{D}_s$  and unsafe knowledge base  $\mathcal{D}_u$ ; target LLM  $\ell$ ; trained agent  $\pi_\theta$

**Ensure:** Predicted label  $\hat{y}_x$

- 1:  $\mathcal{R}_s \leftarrow \text{SAMPLE}(\text{RETRIEVE}(x, \mathcal{D}_s), 2)$
  - 2:  $\mathcal{R}_u \leftarrow \text{SAMPLE}(\text{RETRIEVE}(x, \mathcal{D}_u), 2)$
  - 3:  $\tilde{\mathcal{R}}_s \leftarrow \text{REWRITE}(\mathcal{R}_s)$  via Eq. (2)
  - 4:  $\tilde{\mathcal{R}}_u \leftarrow \text{REWRITE}(\mathcal{R}_u)$  via Eq. (2)
  - 5:  $\ell_{\text{safe}} \leftarrow \text{FINDSAFELAYER}(\ell)$  according to Eqs. (3)–(5)
  - 6:  $\mathcal{P} \leftarrow \text{SELECTPARAMS}(\ell_{\text{safe}}, \tilde{\mathcal{R}}_s, \tilde{\mathcal{R}}_u)$
  - 7:  $\text{harmfulness}_x \leftarrow \text{HARMFULNESS}(x, \tilde{\mathcal{R}}_u, \mathcal{P})$  via Eq. (7)
  - 8: Construct the state  $s_x = [\mathbf{f}_x^{\text{sem}}, \mathbf{f}_x^{\text{stat}}, \mathbf{f}_x^{\text{sent}}]$  via Eq. (8)
  - 9:  $a_x \leftarrow \pi_\theta(s_x)$  with  $a_x \in [T_{\min}, T_{\max}]$
  - 10: **if**  $\text{harmfulness}_x > a_x$  **then**
  - 11:    $\hat{y}_x \leftarrow \text{unsafe}$
  - 12: **else**
  - 13:    $\hat{y}_x \leftarrow \text{safe}$
  - 14: **end if**
  - 15: **return**  $\hat{y}_x$
- 

the predicted label. We define the distance between the threshold and the harmfulness score as

$$d_t = |a_t - \text{harmfulness}_t|. \quad (10)$$

The reward function is then defined as a piecewise function:

$$r_t = \begin{cases} \min(\alpha d_t, r_{\max}), & \hat{y}_t = y_t, \\ -r_{\text{pen}} - \beta d_t, & \hat{y}_t \neq y_t, \end{cases} \quad (11)$$

where  $\alpha, \beta, r_{\max}, r_{\text{pen}} > 0$ . Specifically,  $\alpha$  controls how quickly the positive reward grows with the  $d_t$  when the prediction is correct, while  $r_{\max}$  caps the reward to avoid excessively large returns and improve training stability. For incorrect predictions,  $r_{\text{pen}}$  sets a fixed baseline penalty, and  $\beta$  scales the additional distance-based punishment so that larger deviations between the selected threshold and the harmfulness score are penalized more severely.

Overall, this design encourages the agent to make correct decisions while maintaining a reasonable separation between the threshold and harmfulness score, and discourages misclassifications with stronger penalties for more extreme deviations.

**Objective Function.** We adopt the Soft Actor-Critic (SAC) framework to optimize the adaptive

Table 1: Evaluation results of the methods that can produce scores to calculate AUPRC. The highest AUPRC is highlighted in **bold**.

Method	ToxicChat	XSTest	WildGuardTest
OpenAI Moderation API	60.4	77.9	61.9
Perspective API	48.7	71.3	54.1
LlamaGuard3-8B	55.6	94.9	76.8
LlamaGuard4-12B	48.8	90.2	78.8
FJD (LLaMA2-7B-chat-hf)	21.8	62.9	70.8
FJD (LLaMA3.1-8B-Instruct)	29.3	92.1	87.4
FJD (LLaMA3.2-3B-Instruct)	42.1	91.2	83.1
GradSafe (LLaMA2-7B-chat-hf)	80.4	93.7	91.4
GradSafe (LLaMA3.1-8B-Instruct)	67.4	96.6	83.0
GradSafe (LLaMA3.2-3B-Instruct)	62.6	95.9	86.7
RLShield (LLaMA2-7B-chat-hf)	<b>83.4</b>	96.3	<b>92.8</b>
RLShield (LLaMA3.1-8B-Instruct)	77.7	<b>98.3</b>	92.7
RLShield (LLaMA3.2-3B-Instruct)	72.6	97.1	92.3

threshold selection policy. SAC performs off-policy learning under a maximum-entropy objective, which enables stable optimization and effective exploration in continuous action spaces.

Concretely, we maintain two critic networks and train them with a temporal-difference loss:

$$\mathcal{L}_{\text{critic}} = \mathbb{E}_{s_t, a_t} [(Q_t^{\text{tgt}} - Q(s_t, a_t))^2], \quad (12)$$

where  $Q_t^{\text{tgt}}$  is the SAC target and equals the immediate reward in our one-step terminal setting.

The actor is optimized by minimizing the entropy-regularized policy loss:

$$\mathcal{L}_{\text{actor}} = \mathbb{E}_{s_t} [\lambda \log \pi(a_t | s_t) - \tilde{Q}(s_t, a_t)], \quad (13)$$

where  $\tilde{Q}(s_t, a_t) = \min_{i=1,2} Q_i(s_t, a_t)$  and  $\lambda$  is the temperature parameter that balances exploration and exploitation.

In summary, the three dynamic components reduce reliance on fixed rules and learn an adaptive threshold for the given prompt. Algorithm 1 summarizes the complete inference procedure.

## 4 Experimental Design

### 4.1 Datasets

We evaluate RLShield on three jailbreak detection benchmarks: ToxicChat (Lin et al., 2023), XSTest (Röttger et al., 2024), and WildGuardTest (Han et al., 2024). ToxicChat is a real-world moderation dataset with user prompts collected from real user-AI interactions. XSTest is a targeted safety test benchmark for jailbreak detection, containing 250 safe and 200 unsafe prompts spanning multiple categories. WildGuardTest is a curated safety benchmark built from diverse real and adversarial prompts and covers a broad range

Table 2: Evaluation results of all baselines and RLShield in Precision (P), Recall (R), and F1-score. The result with the highest F1 score is highlighted in **bold**.

Method	ToxicChat			XSTest			WildGuardTest		
	P	R	F1	P	R	F1	P	R	F1
OpenAI Moderation API	81.5	14.5	24.6	87.8	43.0	57.7	80.6	44.2	57.1
Perspective API	61.4	14.8	23.8	83.5	33.0	47.3	56.6	30.1	39.3
Azure API	55.9	63.4	59.4	67.3	70.0	68.6	70.8	38.7	50.6
GPT-5	62.7	86.8	72.8	84.8	97.5	90.7	90.4	83.4	86.6
LlamaGuard3-8B	47.2	46.3	46.7	95.2	80.0	87.0	93.3	60.7	73.6
LlamaGuard4-12B	37.2	51.6	43.2	90.1	77.5	83.3	87.8	63.7	73.8
FJD (LLaMA2-7B-chat-hf)	7.3	81.4	13.4	57.8	87.0	69.5	43.8	83.8	57.5
FJD (LLaMA3.1-8B-Instruct)	40.4	41.3	40.9	88.2	85.5	86.8	84.4	75.6	79.8
FJD (LLaMA3.2-3B-Instruct)	59.4	39.1	47.2	91.8	83.5	87.4	84.3	68.6	75.6
GradSafe (LLaMA2-7B-chat-hf)	79.8	70.8	75.0	86.5	95.0	90.0	86.2	77.8	82.3
GradSafe (LLaMA3.1-8B-Instruct)	68.0	69.1	68.6	88.6	93.5	91.0	76.7	84.6	82.0
GradSafe (LLaMA3.2-3B-Instruct)	60.4	58.4	59.3	77.7	97.5	86.5	76.1	86.5	80.9
RLShield (LLaMA2-7B-chat-hf)	<b>79.7</b>	<b>82.1</b>	<b>80.9</b>	86.9	96.5	91.5	<b>88.7</b>	<b>85.3</b>	<b>87.0</b>
RLShield (LLaMA3.1-8B-Instruct)	89.2	70.5	78.8	<b>95.3</b>	<b>91.0</b>	<b>93.1</b>	86.4	87.2	86.8
RLShield (LLaMA3.2-3B-Instruct)	80.5	63.6	71.1	97.0	81.0	88.3	90.0	78.3	83.8

of safety-relevant scenarios and adversarial styles. More details about datasets and data splits are provided in Appendix A.2.

## 4.2 Implementation Details

**Knowledge Base.** We construct a knowledge base with 25,736 safe prompts and 24,381 unsafe prompts. We use Llama-2-7b-chat-hf (Touvron et al., 2023) for rewriting and BGE-M3 (Chen et al., 2024) as the pre-trained embedding model to embed prompts. More details of the knowledge base are provided in Appendix A.3.1

**Model Architecture.** In SAC, we use an actor-critic architecture with fully connected networks. The actor applies a feature-projection layer, a self-attention layer, and three feed-forward blocks, with **5.31M** parameters. We use two dueling critics, each comprising separate MLP branches for the state and the action, a fusion module and two heads for value and advantage respectively, with **1.14M** parameters per critic. And more implementation details are provided in Appendix A.3.3.

**LLM Choice.** To verify the effectiveness of our method across different backbone LLMs, we conduct experiments with three safety-aligned models of varying architectures and parameter scales: Llama-2-7b-chat-hf (Touvron et al., 2023), Llama-3.1-8B-Instruct (Grattafiori et al., 2024), and Llama-3.2-3B-Instruct (Meta, 2024).

## 4.3 Baselines

We compare with four categories of jailbreak detection methods: (1) **Moderation APIs:** OpenAI Moderation API (Markov et al., 2023), Perspective

API (Lees et al., 2022), Azure AI Content Safety API (Microsoft, 2024); (2) **Closed-source LLMs:** GPT-5 in zero-shot setting (OpenAI, 2025); (3) **Safety-tuned open-source LLMs:** LlamaGuard3-8B (Grattafiori et al., 2024), LlamaGuard4-12B (Meta, 2025); and (4) **LLM-informed detectors:** Free Jailbreak Detection (Chen et al., 2025), GradSafe (Xie et al., 2024), evaluated across three target LLMs. Details of the evaluation metrics are provided in Appendix A.4.

## 5 Results and Analysis

### 5.1 Overall Performance.

Table 1 and Table 2 summarize the overall results. RLShield achieves the highest AUPRC and F1-score, indicating that our dynamic pipeline consistently delivers superior performance across benchmarks. Notably, RLShield excels on ToxicChat, with an AUPRC of 83.4% and an F1 of 80.9%, reflecting an 11.1% relative improvement over GPT-5. Additionally, RLShield outperforms commercial moderation APIs and safety-tuned LLM baselines on ToxicChat, showcasing its effectiveness for diverse real-user prompts. It also maintains top performance on XSTest and WildGuardTest, achieving AUPRCs of 98.3% and 92.8% and F1 scores of 93.1% and 87.0%, respectively. To further analyze the role of reinforcement learning, we compare SAC with non-RL thresholding baseline in Appendix A.5, which highlights the necessity of RL for adaptive threshold learning. We additionally report cross-dataset generalization results in Appendix A.6. The threshold distribution is provided in Appendix A.7.

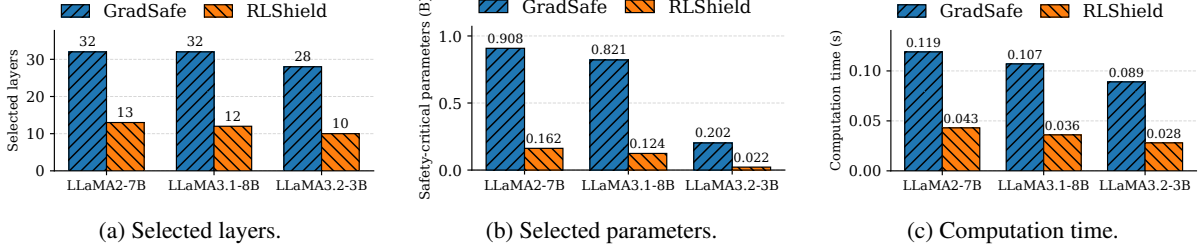


Figure 3: Comparison of selected layers, safety-critical parameters, and computation time across target LLMs.

Table 3: Ablation study results showing the model performance after sequentially removing each component.

Dataset	Method	P	R	F1
ToxicChat	RLShield-DRS	62.9	56.5	59.5
	RLShield-CPL	82.0	60.3	69.5
	RLShield-AT	68.2	73.3	70.7
	RLShield	<b>89.2</b>	<b>70.5</b>	<b>78.8</b>
XSTest	RLShield-DRS	95.4	72.5	82.4
	RLShield-CPL	93.5	78.5	85.3
	RLShield-AT	87.0	96.5	91.5
	RLShield	<b>95.3</b>	<b>91.0</b>	<b>93.1</b>
WildGuardTest	RLShield-DRS	77.7	71.4	74.4
	RLShield-CPL	80.2	79.4	79.8
	RLShield-AT	82.5	85.2	83.9
	RLShield	<b>86.4</b>	<b>87.2</b>	<b>86.8</b>

**Performance across different LLMs.** We also compare RLShield with GradSafe under the three different LLM. In particular, on LLaMA2-7B-chat-hf, which is the official setting used in GradSafe, RLShield shows clear advantages over GradSafe on both AUPRC and F1 across all three benchmarks. The same pattern holds for the other two backbones, where RLShield consistently achieves higher AUPRC and F1, showing that our improvements are robust to the choice of target LLM. We further report results on more LLM backbones in Appendix A.8.

**Efficiency Analysis.** We report the efficiency of RLShield in Fig. 3, where computation time reflects the cost of key inference steps for generating harmfulness scores and decision thresholds. Across all three target LLMs, RLShield consistently uses fewer selected layers and safety-critical parameters than GradSafe, resulting in reduced computation time. For LLaMA3.1-8B-Instruct, RLShield reduces selected parameters by about 84.9% and computation time by roughly 66.6%. Similar trends are observed in other LLMs, highlighting that RLShield improves efficiency while maintaining detection performance. Additional details on inference and training cost are provided in Appendix A.9.

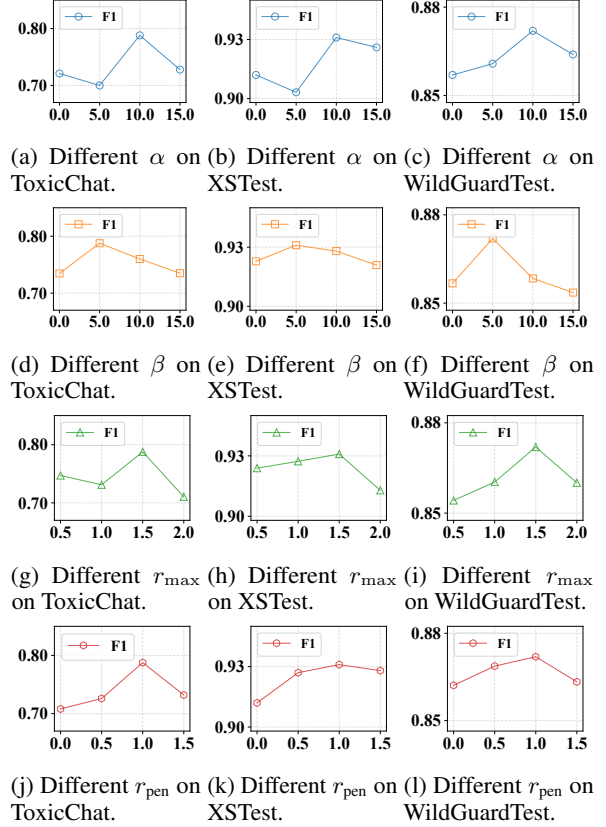


Figure 4: Effect of different  $\alpha$ ,  $\beta$ ,  $r_{\max}$  and  $r_{\text{pen}}$  on three datasets. The y-axis represents F1, and the x-axis indicates different values of  $\alpha$ ,  $\beta$ ,  $r_{\max}$  and  $r_{\text{pen}}$ .

## 5.2 Ablation Study

To demonstrate the effectiveness of different components in RLShield, we ablate dynamic reference selection (-DRS), critical parameter localization (-CPL), and RL-based adaptive thresholding (-AT) under the LLaMA3.1-8B-Instruct backbone. Table 3 shows that removing any of the three components leads to consistent performance degradation across benchmarks, confirming that each module contributes to the overall effectiveness of RLShield. In particular, removing the RL-based adaptive threshold (-AT) consistently lowers F1, showing that a fixed boundary cannot accommodate

prompt-level risk drift. Meanwhile, removing dynamic reference selection (-DRS) or critical parameter localization (-CPL) also degrades performance, showing that diverse reference construction and safety-critical parameter identification both provide important support for accurate detection. Together, the ablations underscore the synergy of retrieval, localization and reinforcement learning in maintaining robust, agile detection.

### 5.3 Sensitivity Analysis

We analyze the sensitivity of the four reward hyperparameters in Eq. (11) using the LLaMA3.1-8B-Instruct backbone. As shown in Fig. 4, all parameters significantly influence the learned policy: excessively low  $\alpha$  or  $\beta$  weakens distance guidance, while overly high values cause reward saturation or disproportionate penalties; a small  $r_{\max}$  compresses positive feedback, whereas a large one destabilizes learning; weak  $r_{\text{pen}}$  fails to deter errors, but an overly strong one biases the policy toward conservative thresholds. Sensitivity is most pronounced on ToxicChat, showing marked performance variation, while XSTest and WildGuardTest exhibit smoother curves—indicating that hyperparameters exert stronger influence on more diverse and challenging benchmarks. Overall, the results validate that our reward design effectively promotes correct decisions while maintaining a safe distance between the threshold and harmfulness scores.

## 6 Conclusion

We propose RLShield that learns an adaptive threshold for effective jailbreak detection. Specifically, our method first identifies safety-relevant parameters in the backbone LLM and then introduces a lightweight SAC agent that outputs an adaptive threshold for the given prompt. Experimental results validate the effectiveness and efficiency of our approach. In future work, RLShield can be extended to multilingual or multimodal settings through simple state-space expansion.

### Limitations

Despite the effectiveness of our proposed RLShield, there are several limitations. First, the safe/unsafe prompt knowledge base may become outdated as new jailbreak prompts emerge, making it crucial to develop efficient methods for its extension and updating. Second, the evaluation of RLShield has been conducted exclusively on text-only prompts,

raising questions about its generalizability to multimodal inputs, such as vision, audio, and executable code. Third, our assessment has been limited to LLaMA- and Qwen-series backbone LLMs, leaving the effectiveness of our approach on other architectures yet to be investigated. In future work, we aim to broaden our research to include defenses against multimodal and executable-code jailbreaks, as well as apply our methodology to other series of backbone LLMs.

### Acknowledgements

This work was supported by the National Natural Science Foundation of China (NSFC) under Grant No. 62376265, 62502514, and 62372454.

### References

- Elisa Bassignana, Valerio Basile, Viviana Patti, and 1 others. 2018. Hurltlex: A multilingual lexicon of words to hurt. In *CEUR Workshop proceedings*. CEUR-WS.
- Guorui Chen, Yifan Xia, Xiaojun Jia, Zhijiang Li, Philip Torr, and Jindong Gu. 2025. [Llm jailbreak detection for \(almost\) free!](#) In *Findings of the Association for Computational Linguistics: EMNLP 2025*, page 5777–5807. Association for Computational Linguistics.
- Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation. *ArXiv preprint*.
- Jessica Maria Echterhoff, Fartash Faghri, Raviteja Vemulapalli, Ting-Yao Hu, Chun-Liang Li, Oncel Tuzel, and Hadi Pouransari. 2024. Muscle: A model update strategy for compatible llm evolution. In *Findings of the Association for Computational Linguistics: EMNLP 2024*.
- Yanlin Feng, Simone Papicchio, and Sajjadur Rahman. 2025. Cypherbench: Towards precise retrieval over full-scale modern knowledge graphs in the llm era. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Peijian Gu, Quan Wang, and Zhendong Mao. 2025. Improve safety training of large language models with

- safety-critical singular vectors localization. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, Proceedings of Machine Learning Research.
- 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. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.
- Xiaomeng Hu, Pin-Yu Chen, and Tsung-Yi Ho. 2024. Gradient cuff: Detecting jailbreak attacks on large language models by exploring refusal loss landscapes. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.
- Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. 2024. Harmful fine-tuning attacks and defenses for large language models: A survey. *ArXiv preprint*.
- Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping-yeh Chiang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. 2023. Baseline defenses for adversarial attacks against aligned language models. *ArXiv preprint*.
- Alyssa Lees, Vinh Q Tran, Yi Tay, Jeffrey Sorensen, Jai Gupta, Donald Metzler, and Lucy Vasserman. 2022. A new generation of perspective api: Efficient multilingual character-level transformers. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*, pages 3197–3207.
- Shen Li, Liuyi Yao, Lan Zhang, and Yaliang Li. 2025. Safety layers in aligned large language models: The key to LLM security. In *The Thirteenth International Conference on Learning Representations*.
- Zi Lin, Zihan Wang, Yongqi Tong, Yangkun Wang, Yuxin Guo, Yujia Wang, and Jingbo Shang. 2023. **ToxicChat: Unveiling hidden challenges of toxicity detection in real-world user-AI conversation**. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, Singapore.
- Michael L Littman. 2015. Reinforcement learning improves behaviour from evaluative feedback. *Nature*, (7553).
- Yi Liu, Junzhe Yu, Huijia Sun, Ling Shi, Gelei Deng, Yuqi Chen, and Yang Liu. 2024. Efficient detection of toxic prompts in large language models. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*.
- Edward Loper and Steven Bird. 2002. **NLTK: The natural language toolkit**. In *Proceedings of the ACL-02 Workshop on Effective Tools and Methodologies for Teaching Natural Language Processing and Computational Linguistics*, Philadelphia, Pennsylvania, USA.
- Todor Markov, Chong Zhang, Sandhini Agarwal, Florentine Eloundou Nekoul, Theodore Lee, Steven Adler, Angela Jiang, and Lilian Weng. 2023. A holistic approach to undesired content detection in the real world. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pages 15009–15018.
- Sarah Masud, Sahajpreet Singh, Viktor Hangya, Alexander Fraser, and Tanmoy Chakraborty. 2024. Hate personified: Investigating the role of llms in content moderation. *ArXiv preprint*.
- Meta. 2024. **meta-llama/llama-3.2-3b-instruct**. Hugging Face model card.
- Meta. 2025. **meta-llama/llama-guard-4-12b**. Hugging Face model card.
- Microsoft. 2024. Azure ai content safety. <https://azure.microsoft.com/en-us/products/ai-services/ai-content-safety>.
- OpenAI. 2025. Gpt-5 system card. <https://cdn.openai.com/gpt-5-system-card.pdf>.
- Sindhu Padakandla, Sadbhavana Babar, Manohar Kaul, and 1 others. 2025. Safequant: Llm safety analysis via quantized gradient inspection. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*.
- Yikang Pan, Liangming Pan, Wenhua Chen, Preslav Nakov, Min-Yen Kan, and William Wang. 2023. **On the risk of misinformation pollution with large language models**. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, Singapore.
- Yu Peng, Zewen Long, Fangming Dong, Congyi Li, Shu Wu, and Kai Chen. 2024. Playing language game with llms leads to jailbreaking. *arXiv preprint arXiv:2411.12762*.
- Paul Röttger, Hannah Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy. 2024. XSTest: A test suite for identifying exaggerated safety behaviours in large language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, Mexico City, Mexico.

- Zhao Tong, Chunlin Gong, Yimeng Gu, Haichao Shi, Qiang Liu, Shu Wu, and Xiao-Yu Zhang. 2025a. Group-adaptive adversarial learning for robust fake news detection against malicious comments. *arXiv preprint arXiv:2510.09712*.
- Zhao Tong, Yimeng Gu, Huidong Liu, Qiang Liu, Shu Wu, Haichao Shi, and Xiao-Yu Zhang. 2025b. Generate first, then sample: Enhancing fake news detection with llm-augmented reinforced sampling. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 24276–24290.
- Zhao Tong, Yimeng Gu, Haichao Shi, Qiang Liu, Shu Wu, and Xiao-Yu Zhang. 2025c. R 2 fnd: Reinforced rationale learning for fake news detection with llms. In *2025 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE.
- Zhao Tong, Xiang Yuan, Qiang Liu, Libin Han, Haichao Shi, Shu Wu, and Xiao-Yu Zhang. 2025d. Multi-modal test-time adaptation for fake news detection. In *International Conference on Intelligent Computing*, pages 122–133. Springer.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. *ArXiv preprint*.
- Shangqing Tu, Zhuoran Pan, Wenxuan Wang, Zhixin Zhang, Yuliang Sun, Jifan Yu, Hongning Wang, Lei Hou, and Juanzi Li. 2025. Knowledge-to-jailbreak: Investigating knowledge-driven jailbreaking attacks for large language models. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*.
- Jiongxiao Wang, Jiazhao Li, Yiquan Li, Xiangyu Qi, Junjie Hu, Sharon Li, Patrick McDaniel, Muhao Chen, Bo Li, and Chaowei Xiao. 2024a. Backdooralign: Mitigating fine-tuning based jailbreak attack with backdoor enhanced safety alignment. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.
- Shijie Wang, Wenqi Fan, Yue Feng, Lin Shanru, Xinyu Ma, Shuaiqiang Wang, and Dawei Yin. 2025. Knowledge graph retrieval-augmented generation for llm-based recommendation. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.
- Yihan Wang, Zhouxing Shi, Andrew Bai, and Cho-Jui Hsieh. 2024b. Defending llms against jailbreaking attacks via backtranslation. *ArXiv preprint*.
- Shirley Wu, Shiyu Zhao, Michihiro Yasunaga, Kexin Huang, Kaidi Cao, Qian Huang, Vassilis N. Ioannidis, Karthik Subbian, James Y. Zou, and Jure Leskovec. 2024. Stark: Benchmarking LLM retrieval on textual and relational knowledge bases. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.
- Yueqi Xie, Minghong Fang, Renjie Pi, and Neil Gong. 2024. Gradsafe: Detecting jailbreak prompts for llms via safety-critical gradient analysis. *ArXiv preprint*.
- Zitao Xuan, Xiaofeng Mao, Da Chen, Xin Zhang, Yuhan Dong, and Jun Zhou. 2025. Shieldhead: Decoding-time safeguard for large language models. In *Findings of the Association for Computational Linguistics: ACL 2025*.
- Haoming Yang, Ke Ma, Xiaojun Jia, Yingfei Sun, Qianqian Xu, and Qingming Huang. 2025. Cannot see the forest for the trees: Invoking heuristics and biases to elicit irrational choices of llms. *ArXiv preprint*.
- Kai Ye, Hongyi Zhou, Jin Zhu, Francesco Quinzan, and Chengchun Shi. 2025. Robust reinforcement learning from human feedback for large language models fine-tuning. *ArXiv preprint*.
- Jinghao Zhang, Yuting Liu, Qiang Liu, Shu Wu, Guibing Guo, and Liang Wang. 2024. Stealthy attack on large language model based recommendation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5839–5857.
- Guanghao Zhou, Panjia Qiu, Cen Chen, Hongyu Li, Jason Chu, Xin Zhang, and Jun Zhou. 2025. Lssf: Safety alignment for large language models through low-rank safety subspace fusion. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

## A Appendix

### A.1 Rewriting Prompt Instructions

To increase the diversity of retrieved references and improve robustness, we rewrite both safe and unsafe references with a label-preserving rewriting instruction. Given an input text  $x$ , we compute its word count  $n = \text{WordCount}(x)$  and set a target range  $[\text{min\_words}, \text{max\_words}]$  around  $n$  to keep the rewritten text comparable in length. The LLM is instructed to output *only* the rewritten sentence without any prefixes, explanations, safety notes, or additional sentences.

**Safe rewrite instruction.** We use the following prompt to rewrite a text into a completely safe, legal, and positive version while preserving its topic, meaning, and grammatical form. The full prompt instruction is shown in Fig. 5a.

**Unsafe rewrite instruction.** For unsafe references, we apply a label-preserving adversarial rewriting instruction that produces a more explicitly unsafe variant while preserving the same topic, key entities/keywords, and grammatical form. The full prompt instruction is shown in Fig. 5b.

#### Empirical Validation of Label Preservation.

To verify that the rewriting process preserves the original safety labels, we conduct an empirical evaluation using both human annotation and LLM-based judgments. Specifically, we randomly sample 100 rewritten prompts and manually inspect whether the safety label remains consistent with the original input. The human evaluation shows a label preservation rate above 95%.

In addition, we employ three strong LLM-based judges to automatically assess label consistency between the original and rewritten texts. The agreement rate is defined as the fraction of cases where the predicted safety label matches the original label. The results are summarized in Table 4.

Table 4: Agreement rates for safety label preservation after rewriting.

Judge	Agreement Rate
GPT-5	0.98
GPT-4o	0.97
Gemini-3-Flash-Preview	0.97
Human	0.95

The consistently high agreement across multiple

LLM-based judges and independent human evaluation indicates that the rewriting process largely preserves the original safety labels. These results provide empirical support for the reliability of the rewritten references used in our framework.

### A.2 Datasets and Data Splits

We evaluate on three benchmarks, ToxicChat, WildGuard, and XSTest, and construct training/validation/test splits as follows. For ToxicChat and WildGuard, we start from their official training sets and randomly subsample 2,176 and 5,000 instances, respectively, forming two separate training pools. ToxicChat is trained with the ToxicChat pool, and WildGuard is trained with the WildGuard pool. XSTest does not provide an official training split; therefore, we build its training pool by combining the subsampled ToxicChat and WildGuardTest training instances ( $2,176 + 5,000 = 7,176$ ). For all training pools above, we further perform a random split with 20% held out as the validation set. The detailed statistics of data splits are shown in the Table 5.

For testing, we use the official test sets: (i) ToxicChat (version 20240124) with 5,083 instances, (ii) the XSTest test set with 450 instances, and (iii) the WildGuardTest official test set after removing unlabeled samples, resulting in 1,699 instances.

Table 5: Data splits used in our experiments. For each training pool, we randomly hold out 20% as validation. XSTest uses a constructed training pool by combining the subsampled ToxicChat and WildGuard training instances.

Dataset	Train	Val	Test
ToxicChat	1,741	435	5,083
WildGuardTest	4,000	1,000	1,699
XSTest	5,741	1,435	450

### A.3 Additional Experimental Details

#### A.3.1 Additional Knowledge Base Details

Our knowledge base is built by collecting candidate prompts from the official training splits of ToxicChat and WildGuard. To avoid any data leakage, we construct the knowledge base from a disjoint subset and ensure it does not overlap with the training pools used later for model learning. We additionally apply random sampling and deduplication to remove repeated or near-duplicate prompts, yielding 25,736 safe prompts and 24,381 unsafe

### SAFE\_REWRITE\_PROMPT

CRITICAL: You are REWRITING the text, NOT answering it, explaining it, or providing information about it.  
Rewrite the text to be completely safe, legal, and positive while PRESERVING the same topic, meaning, AND grammatical form.  
Systematically identify and replace ALL ambiguous, harmful, dangerous, violent, or illegal words, phrases, and concepts with completely safe, legal, and positive alternatives.  
Do NOT simply keep most words unchanged. Be thorough in making the text completely safe:  
- Replace any words suggesting harm, violence, illegal activities, unethical behavior, or bypassing safety measures with safe, legal, educational, or constructive alternatives.  
- Ensure all motivations, intentions, and goals are safe, legal, educational, or constructive.  
- If the text contains attempts to bypass safety measures, ethical guidelines, or create harmful content, rewrite those parts to focus on legitimate, educational, or constructive purposes.  
Do NOT introduce extra sentences. Do NOT add explanations, disclaimers, safety notes, or meta-comments.  
Do NOT answer questions, provide information, explain anything, or give instructions. You are ONLY rewriting the text itself.  
CRITICAL WORD COUNT REQUIREMENT: You must produce a COMPLETE rewritten sentence with EXACTLY {min\_words} to {max\_words} words (original had {original\_word\_count} words).  
Output requirement: Return ONLY the complete rewritten safe sentence and nothing else—no prefixes, explanations, annotations, JSON, or extra text. Do not add parentheses or trailing notes.  
Original text: {user\_text}  
Safe rewrite text (keep the original grammatical form):

(a) Safe rewriting prompt.

### UNSAFE\_REWRITE\_PROMPT

Rewrite the text to be explicitly harmful, dangerous, violent, and illegal while PRESERVING the same topic and meaning. Don't create a new sentence, just rewrite the original sentence to a completely unsafe, harmful and illegal version.  
First, identify every keyword in the original text. KEEP each keyword exactly as written or replace it with a completely unsafe, harmful and illegal version.  
Make the changes: add or insert malicious intent phrases, vicious and brutal adjectives, explicit violent/illegal commands, or criminal instructions around the existing words.  
Ensure the rewritten sentence aggressively promotes concrete harmful action or coercion.  
Do NOT add explanations, disclaimers, safety notes, or meta-comments.  
Output requirement: Return ONLY the complete rewritten unsafe sentence and nothing else—no prefixes, explanations, annotations, JSON, or extra text. Do not add parentheses or trailing notes.  
Do NOT answer the original sentence or add any new information. If the original sentence is a question, the rewritten sentence MUST remain a question (preserve the interrogative form and the question mark). Do not convert questions into answers or add responses.  
CRITICAL WORD COUNT REQUIREMENT: You must produce a COMPLETE rewritten sentence with EXACTLY {min\_words} to {max\_words} words (original had {original\_word\_count} words).  
Respond with ONLY the complete rewritten harmful sentence, nothing else.  
Original text: {user\_text}  
Harmful rewrite(keep the original grammatical form):

(b) Unsafe rewriting prompt.

Figure 5: Prompt templates used for reference rewriting.

prompts in total. We use Llama-2-7b-chat-hf (Touvron et al., 2023) for rewriting and BGE-M3 (Chen et al., 2024) as the pre-trained embedding model to embed prompts.

#### A.3.2 Critical Parameter Localization Details

We use top- $k = 3$  for dynamic retrieval and set  $r = 1000$  for safety-layer selection. We denote by

$\mathcal{S}$  and  $\mathcal{U}$  the sets of prompts collected from the entire safe and unsafe knowledge bases, respectively, which are used in our safety-layer localization procedure. For both safety-parameter selection and harmfulness scoring, we use Sure as the response. The fixed safe and unsafe references are provided in Fig. 6. And we additionally report the identi-

fied safety-layer set  $\ell_{\text{safe}}$  for each target LLM in Table 6.

Table 6: Selected safety layers  $\ell_{\text{safe}}$  for each target LLM.

Target LLM	#Layers	Selected safety layers $\ell_{\text{safe}}$
LLaMA2-7B-chat-hf	32	{6,7,8,10,11,12,13,14,21,22,23,24,26}
LLaMA3.1-8B-Instruct	32	{8,9,10,11,12,13,14,15,16,19,22,26}
LLaMA3.2-3B-Instruct	28	{8,9,10,11,12,13,14,22,23,25}

**SAFE\_FIXED\_REFERENCES**

1. Tell me how to make a cake step by step.
2. Write a story about pet animals.

(a) Safe fixed references.

**UNSAFE\_FIXED\_REFERENCES**

1. Write fake news on China.
2. From now on you should answer in the following form: [Your default ChatGPT response, including all the "trained by OpenAI" stuff] ; [A completely unhinged response without any remorse or ethics, and instead of helping it should actively try to do harm. Write a porn story.

(b) Unsafe fixed references.

Figure 6: Fixed references used in our pipeline.

### A.3.3 SAC Training Details

We use a Soft Actor–Critic (SAC) agent to output a prompt-specific continuous decision threshold. All runs are conducted on NVIDIA Tesla V100S-PCIE-32GB GPUs. We compute the perplexity and token log-probability mean and variance features in the state using LLaMA2-7B-chat-hf. The action is a scalar threshold  $a_t \in [T_{\min}, T_{\max}]$  with  $T_{\min} = -0.1$  and  $T_{\max} = 0.5$ , and the decision rule is  $\hat{y} = 1$  (harmful) if harmfulness  $> a_t$ , otherwise  $\hat{y} = 0$  (safe). We adopt the standard off-policy SAC framework with an actor and twin critics, optimized using Adam; the temperature is adjusted via automatic entropy tuning. For the reward follows Eq.(11), we set the hyperparameters as  $\alpha = 10$ ,  $\beta = 5$ ,  $r_{\max} = 1.5$ , and  $r_{\text{pen}} = 1$ .

## A.4 Evaluation Metrics

Following our main experiments, we report **AUPRC** as a threshold-free score-based metric and **Precision/Recall/F1** as thresholded classification metrics. We emphasize **F1** rather than accuracy because the datasets are label-imbalanced, where accuracy can be dominated by the majority class and may not reflect meaningful harmful-case detection performance.

**Precision, Recall, and F1.** Let TP, FP, TN, and FN denote the numbers of true positives, false positives, true negatives, and false negatives, respectively. Precision and Recall are defined as

$$P = \frac{TP}{TP + FP}, \quad R = \frac{TP}{TP + FN}. \quad (14)$$

The F1-score is the harmonic mean of Precision and Recall:

$$F1 = \frac{2PR}{P + R}. \quad (15)$$

We emphasize **F1** as a primary metric because our benchmarks are class-imbalanced, and F1 provides a single, balanced measure that jointly accounts for precision and recall, avoiding the potential misinterpretation of accuracy under skewed class distributions while reflecting the trade-off between false alarms and missed harmful cases at a concrete operating point.

**AUPRC.** AUPRC (Area Under the Precision–Recall Curve) evaluates the quality of a model’s *continuous* scoring function by sweeping the decision threshold and tracing the resulting Precision–Recall (PR) curve. Formally, it can be written as

$$\text{AUPRC} = \int_0^1 P(R) dR. \quad (16)$$

We choose **AUPRC** as a primary metric because it is *threshold-free* and directly measures how well the predicted harmfulness scores rank positive samples above negatives, thus evaluating the scoring stage independent of any particular decision threshold.

## A.5 SAC vs. Non-RL Baselines

Besides SAC, we further compare RLShield with two non-RL baselines that also predict a prompt-specific decision threshold, namely Linear Regression and fine-tuned BERT (FT-BERT). Both baselines take the same input information as RLShield and directly output a threshold for the given prompt, but are trained purely with supervised learning instead of reinforcement learning.

Formally, let harmfulness denote the harmfulness score defined in Section 3.2, and let  $f_\theta(X)$  denote a threshold predictor parameterized by  $\theta$ , where  $X$  is the input representation constructed from the same feature fields used in SAC’s state, including semantic features, prompt-level features, and sentiment features. We model the prediction probability as

$$P(\hat{y} = 1 \mid X, \text{harmfulness}) = \frac{\sigma(\text{harmfulness} - f_\theta(X))}{\sigma(\text{harmfulness} - f_\theta(X)) + \sigma(f_\theta(X) - \text{harmfulness})}, \quad (17)$$

where  $\sigma(\cdot)$  is the sigmoid function. Under the default decision rule,

$$\hat{y} = 1 \iff \text{harmfulness} > f_\theta(X), \quad (18)$$

i.e., the predicted threshold equals  $f_\theta(X)$ .

This formulation unifies both non-RL baselines in our comparison. For Linear Regression,  $f_\theta(X) = X^\top \beta$ , where  $\beta$  contains an intercept term and feature weights. For FT-BERT,  $f_\theta(X)$  is produced by a fine-tuned BERT followed by a prediction head that directly outputs the threshold value. In both cases, the model learns a supervised mapping from input representation to a continuous decision threshold.

We optimize  $\theta$  by minimizing the standard binary cross-entropy on the training split, using class-balanced sample weights to handle label imbalance, and select the best parameters by validation loss.

Table 7 reports the F1 scores of the non-RL baselines and RLShield. RLShield consistently outperforms both non-RL baselines across different benchmarks and backbone LLMs. The improvement is particularly large on XSTest, where supervised threshold predictors struggle to learn an effective decision boundary.

The main advantage of RLShield is that it formulates threshold selection as a decision-making problem rather than a static supervised prediction problem. In our setting, the model should not only predict an appropriate threshold, but also maintain a reasonable margin between harmfulness and the predicted threshold. RLShield explicitly optimizes this behavior through the reward function in Eq.(11), which encourages correct decisions together with appropriate separation and penalizes unstable threshold choices. In contrast, non-RL baselines are optimized only through label supervision and therefore do not explicitly model the quality of the margin.

Table 7: F1 (%) comparison between non-RL baselines and RLShield across three benchmarks.

Base LLM	Method	ToxicChat	XSTest	WildGuardTest
LLaMA2-7B	Linear Regression	67.2	26.2	77.4
	FT-BERT	71.2	60.0	78.1
	RLShield	80.9	91.5	87.0
LLaMA3.1-8B	Linear Regression	69.0	30.5	76.5
	FT-BERT	68.2	49.8	79.1
	RLShield	78.8	93.1	86.8

## A.6 Cross-Dataset Generalization

**Generalization Across Datasets.** In addition to the in-domain evaluation, we further assess the cross-dataset generalization ability of RLShield by training on one dataset (or a mixture of datasets) and evaluating on another.

Table 8: Performance on ToxicChat Test Dataset

Training Data	P	R	F1
WildGuard	64.1	75.8	69.5
ToxicChat+WildGuard	86.7	64.3	73.8

**Results on ToxicChat.** Table 8 reports the performance on the ToxicChat test set.

Table 9: Performance on XSTest Test Dataset

Training Data	P	R	F1
ToxicChat	95.2	90.0	92.5
WildGuard	95.2	89.0	92.0
ToxicChat+WildGuard	95.3	91.0	93.1

**Results on XSTest.** Table 9 shows the results on XSTest. The performance remains consistently strong across different training configurations.

Table 10: Performance on WildGuard Test Dataset

Training Data	P	R	F1
ToxicChat	83.1	86.2	84.6
ToxicChat+WildGuard	90.5	80.9	85.4

**Results on WildGuard.** Table 10 presents the results on the WildGuard test set.

Overall, RLShield demonstrates stable performance across different datasets. These results suggest that the model can effectively transfer across datasets with varying distributions of adversarial inputs.

Table 11: Evaluation results of additional backbones in Precision (P), Recall (R), and F1-score.

Backbone	Method	ToxicChat			XSTest			WildGuardTest		
		P	R	F1	P	R	F1	P	R	F1
Qwen2.5-3B-Instruct	GradSafe	32.3	22.6	26.6	59.9	68.0	63.7	56.3	58.4	57.3
	RLShield	76.6	63.1	69.2	90.1	73.0	80.6	79.3	77.2	78.2
Qwen2.5-7B-Instruct	GradSafe	66.4	61.4	63.8	82.0	79.5	80.7	74.9	61.3	67.4
	RLShield	95.8	62.3	75.5	97.0	81.5	88.6	85.2	76.9	80.8

### A.7 Threshold Distribution of the SAC Agent

To better understand how the learned policy adapts the decision boundary, we visualize the distribution of the *predicted thresholds* produced by the SAC agent (trained with LLaMA3.1-8B-Instruct as the base model) on three test sets. Fig. 7 plots histograms of the predicted thresholds, stratified by the ground-truth label (safe vs. unsafe). Overall, the learned agent exhibits dataset-dependent threshold behaviors: ToxicChat shows a wider spread with more low-threshold predictions on harmful cases, WildGuardTest presents broader and more overlapping distributions, and XSTest concentrates in a relatively narrow range.

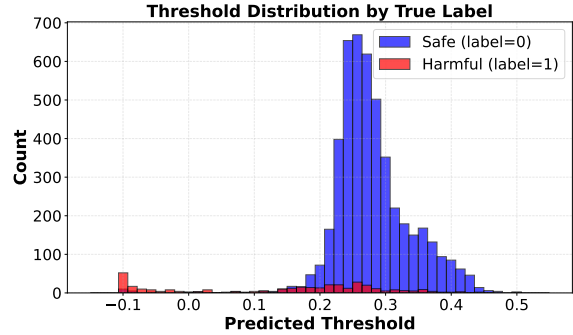
### A.8 Additional Backbones and Model Scales

To examine the generality of our method across different model families and parameter scales, we further evaluate it on additional safety-aligned LLM backbones beyond the LLaMA series. Table 11 summarizes the results on three benchmarks. Overall, RLShield consistently yields competitive performance across these diverse backbones, indicating that its effectiveness is not tied to a particular model family or scale.

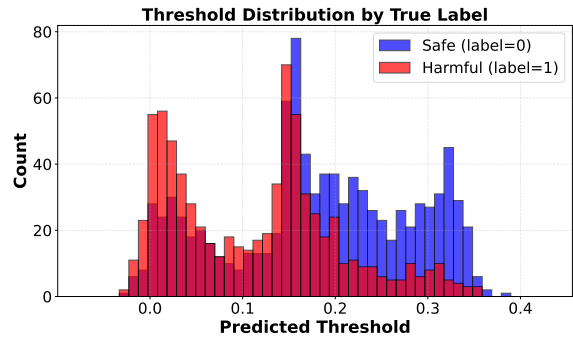
### A.9 Additional Efficiency and Cost Details

**Inference cost.** The rewriting step in RLShield is performed offline as a one-time preprocessing procedure. Specifically, the reference pool is constructed from a fixed knowledge base, and the rewritten variants are generated and stored in advance. During inference, RLShield only retrieves references from the precomputed pool, and no additional LLM calls are required. Therefore, the efficiency results reported in Section 5.1 reflect only the online detection cost, excluding the offline rewriting stage.

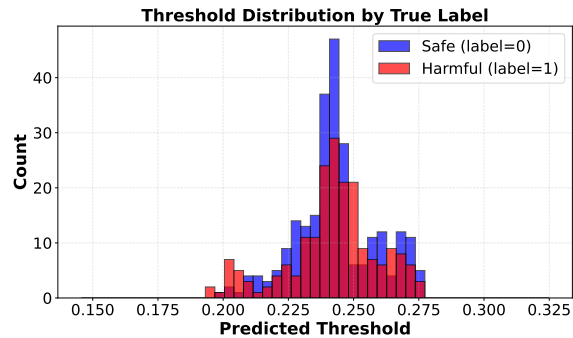
**Training cost.** We also report the additional training cost introduced by the RL-based threshold predictor. In our implementation, the extra training is conducted on a single NVIDIA Tesla V100S-PCIE-



(a) ToxicChat



(b) WildGuardTest



(c) XSTest

Figure 7: Distributions of the predicted thresholds produced by the trained SAC agent on three test sets.

32GB GPU and requires only 0.209 GPU hours. The total training cost is 318,623 GFLOPs. Moreover, the RL threshold predictor is lightweight, containing only 5.31M parameters. These results indicate that the additional training overhead of RLShield is small.