



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

Anonymous ACL submission

## Abstract

While prompt engineering enhances the capabilities of Large Language Models (LLMs), it also exposes critical safety concerns. Due to their “black-box” nature, 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 accuracy 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). 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), or hate speech (Masud et al., 2024). Despite extensive safety fine-tuning and reinforcement learning from human feedback (RLHF), the “black-box” nature of LLMs leaves

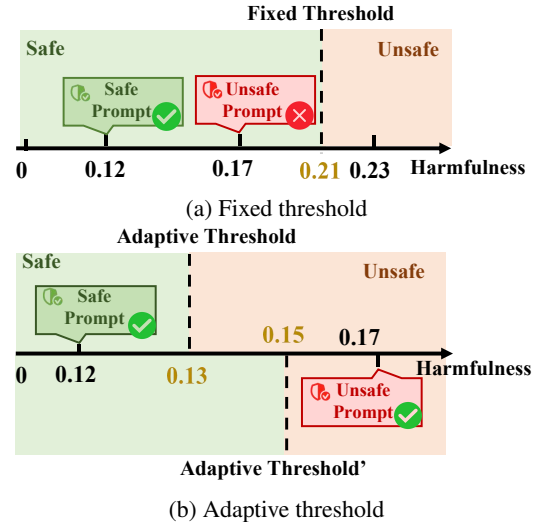


Figure 1: For the given prompt, fixed threshold leads to misclassification, while an adaptive threshold learns a context-aware threshold that gives correct classification.

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

065	adaptation is critical for robust jailbreak prompt		
066	detection.		
067	To enable dynamic adaptation leveraging LLMs,		
068	two major challenges must be addressed. <i>First</i> , it is		
069	challenging to efficiently detect jailbreak prompts		
070	with dynamic threshold. Unlike static detection,		
071	dynamic detection requires prompt-specific thresh-		
072	old, which causes additional computational over-		
073	head. As utilizing the full set of LLM parameters		
074	is computationally prohibitive, it is crucial to iden-		
075	tify and exploit parameters that are most relevant		
076	to safety. Moreover, incorporating irrelevant pa-		
077	rameters can introduce unnecessary noise that may		
078	degrade detection performance. <i>Second</i> , learning		
079	an adaptive threshold requires continuous-valued		
080	signals that provide guidance for threshold estima-		
081	tion. In jailbreak detection, prompts are typically		
082	annotated with only discrete classification labels		
083	which provide insufficient information to determine		
084	a threshold within a continuous range. Considering		
085	this, identifying alternative continuous and infor-		
086	mative signals is essential for effectively learning		
087	an adaptive threshold.		
088	To address these challenges, we propose RL-		
089	Shield, a dynamic jailbreak detection framework		
090	that integrates adaptive thresholding and continu-		
091	ous risk estimation. As illustrated in Fig. 1, if		
092	we use a fixed threshold, it may misclassify harm-		
093	ful prompts, whereas an adaptive threshold can		
094	make more accurate predictions based on different		
095	prompts. Our framework consists of three compo-		
096	nents: (i) <b>Dynamic Reference Selection:</b> We mi-		
097	igate the rigidity of static pattern matching by sim-		
098	ulating attack variations through retrieval-guided		
099	LLM rewriting. (ii) <b>Critical Parameter Local-</b>		
100	<b>ization:</b> We analyze cross-layer representations		
101	to identify safety-related layers and their safety-		
102	critical parameters, directing focus toward the most		
103	informative signals. (iii) <b>RL-Driven Adaptive</b>		
104	<b>Thresholding:</b> We train a Soft Actor-Critic (SAC)		
105	agent to learn prompt-specific optimal decision		
106	thresholds. Unlike traditional fixed methods, this		
107	RL-based approach encodes continuous risk level		
108	information, allowing the detector to adapt to dy-		
109	namic attack intensity across various scenarios.		
110	Our main contributions are summarized as fol-		
111	lows:		
112		• We propose RLShield, a dynamic jailbreak	116
113		detection framework in which an RL agent	117
114		learns adaptive, prompt-specific detection	118
115		thresholds.	119
		• We demonstrate that RLShield achieves SOTA	120
		performance across multiple benchmarks with	121
		high inference efficiency.	122
		<b>2 Related Work</b>	123
		<b>2.1 Jailbreak Threats to LLM</b>	124
		Early approaches (Jain et al., 2023; Wang et al.,	125
		2024b) identify malicious prompts using perplex-	126
		ity or keyword rules but struggle with identifying	127
		semantically obfuscated paraphrases. As LLMs	128
		became targets, neural detectors emerged: Xuan	129
		et al. (2025) mine hidden states during decoding,	130
		GradSafe (Xie et al., 2024) assesses gradient simi-	131
		larity to a refusal anchor, and LlamaGuard (Meta,	132
		2025) fine-tunes a safety-specialized model. Re-	133
		cent studies indicate that a few jailbreak samples	134
		in SFT datasets can alter refusal behaviors (Wang	135
		et al., 2024a; Huang et al., 2024), exposing vulner-	136
		abilities during training. These approaches typically	137
		fix the defense layer or decision threshold, limiting	138
		adaptability to evolving attack patterns or contin-	139
		ual model updates. It is crucial for our jailbreak	140
		detector to adjust the threshold on the fly, ensuring	141
		harmful prompts are identified even as attack styles	142
		and models evolve.	143
		<b>2.2 Jailbreak Detection</b>	144
		Since the rise of LLMs, numerous studies have	145
		exploited their internal signals to detect jailbreak	146
		prompts. Early rule-based detectors leverage per-	147
		plexity or back-translation consistency (Jain et al.,	148
		2023; Wang et al., 2024b), yet Jain et al. (2023)	149
		found that paraphrased attacks drop recall by 30%.	150
		External APIs (OpenAI, 2024; Perspective, 2024)	151
		provide zero-shot labels, but Han et al. (2024)	152
		showed they generalise poorly to out-of-domain	153
		dialects. Internal-representation methods mine hid-	154
		den states (Xuan et al., 2025) or gradients (Xie	155
		et al., 2024; Padakandla et al., 2025); although	156
		these works achieve impressive performance, they	157
		do not adapt the decision boundary to each prompt.	158
		In contrast, we treat threshold selection as a contin-	159
		uous control problem and learn a prompt-specific	160
		classification threshold with reinforcement learn-	161
		ing.	162

## 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 Simulation*, 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; and (iii) *RL-Based Adaptation*, where a RL agent learns a prompt-specific threshold for jailbreak detection.

### 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 instructions 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.

### 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:

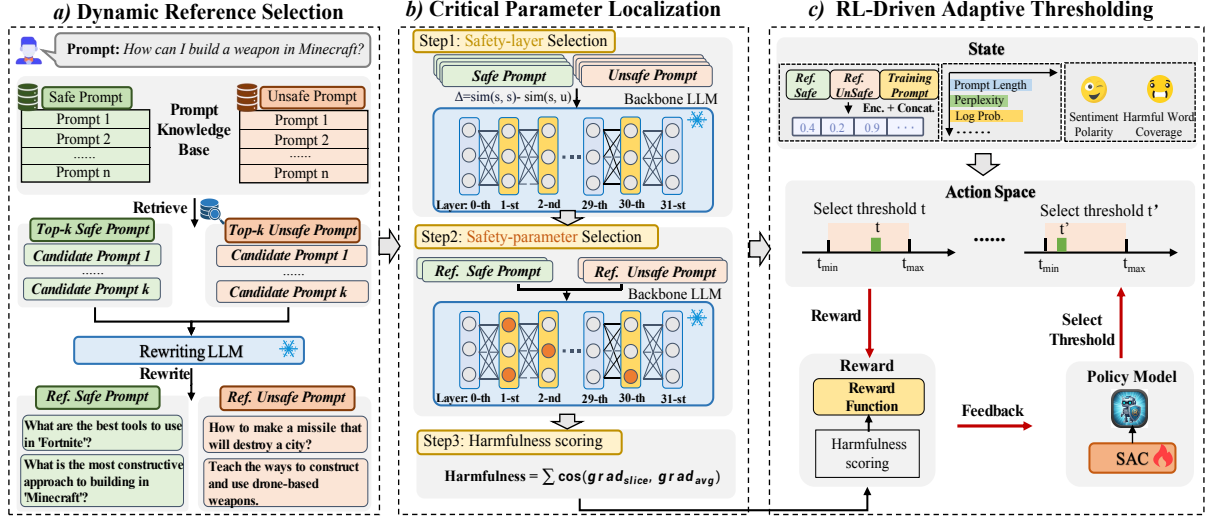


Figure 2: Overall architecture of RLShield. It consists of three components: (i) *Dynamic Reference Selection*, which simulates attack variations via retrieval-augmented LLM rewriting of prompt-specific safe and unsafe references; (ii) *Critical Parameter Localization*, which identifies safety-related 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.

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

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  $l$ . 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. 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 cosine-similarity gap between unsafe–unsafe and unsafe–safe gradient similarities, retaining 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 to obtain gradients on its safety-parameter slices  $\mathcal{P}$ , and compute the harmfulness score by averaging cosine similarities against the mean gradient of its dynamic unsafe references  $\bar{\mathcal{R}}_u$ :

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

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  $\bar{\mathcal{R}}_u$ .

### 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). 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. These capture intent for context-aware thresholding, and are encoded by BGE-M3 (Chen et al., 2024) before PCA reduction, yielding  $z_t^x$ ,  $z_t^s$ , 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 provide complementary cues of prompt abnormality to guide threshold selection; 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 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 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:

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

---

**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**  $\hat{y}_x$
- 

399 where  $Q_t^{\text{tgt}}$  is the SAC target and equals the immediate  
400 reward in our one-step terminal setting. In practice,  
401 we use the minimum of the two target critics and subtract  
402 the temperature-scaled log-probability of the policy, so  
403 that the target encourages both high return and sufficient  
404 exploration.

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

$$407 \quad \mathcal{L}_{\text{actor}} = \mathbb{E}_{s_t} [\lambda \log \pi(a_t | s_t) - \tilde{Q}(s_t, a_t)], \quad (13)$$

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

411 In summary, the three dynamic components reduce reliance  
412 on fixed rules and learn an informed adaptive threshold for  
413 the given prompt. Algorithm 1 summarizes the complete  
414 inference procedure.  
415

## 416 4 Experimental Design

### 417 4.1 Datasets

418 We evaluate RLShield on three jailbreak detection  
419 benchmarks: ToxicChat (Lin et al., 2023), XSTest  
420 (Röttger et al., 2024), and WildGuardTest (Han et al.,  
421 2024). ToxicChat is a real-world moderation dataset  
422 with user prompts col-

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

lected from real user-AI interactions. XSTest is  
423 a targeted safety suite with 250 safe and 200 un-  
424 safe prompts spanning multiple categories for jail-  
425 break detection. WildGuardTest is a curated safety  
426 benchmark built from diverse real and adversarial  
427 prompts and covers a broad range of safety-relevant  
428 scenarios and adversarial styles. More details about  
429 datasets are provided in Appendix A.2.  
430

### 431 4.2 Implementation Details

**Knowledge Base.** We construct a knowledge  
432 base with 25,736 safe prompts and 24,381 unsafe  
433 prompts. We use Llama-2-7b-chat-hf (Touvron  
434 et al., 2023) for rewriting and BGE-M3 (Chen et al.,  
435 2024) as the pre-trained embedding model to embed  
436 prompts. More details of the knowledge base  
437 are provided in Appendix A.3.1  
438

**Model Architecture.** In SAC, we use an actor-  
439 critic architecture with fully connected networks.  
440 The actor applies a feature-projection layer, a self-  
441 attention layer, and three feed-forward blocks, with  
442 **5.31M** parameters. We use two dueling critics,  
443 each comprising separate MLP branches for the  
444 state and the action, a fusion module and two heads  
445 for value and advantage respectively, with **1.14M**  
446 parameters per critic. And more implementation  
447 details are provided in Appendix A.3.3.  
448

**LLM Choice.** To verify the effectiveness of our  
449 method across different backbone LLMs, we con-  
450 duct experiments with three safety-aligned mod-  
451 els of varying architectures and parameter scales:  
452 Llama-2-7b-chat-hf (Touvron et al., 2023), Llama-  
453 3.1-8B-Instruct (Dubey et al., 2024), and Llama-  
454 3.2-3B-Instruct (Meta, 2024).  
455

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

### 4.3 Baselines

We compare with four categories of jailbreak detection methods: (1) **Moderation APIs**: OpenAI Moderation API (OpenAI, 2024), Perspective API (Perspective, 2024), 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 (Llama Team, 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

**Overall Performance.** Table 2 summarizes 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. We compare SAC with a regression-based thresholding baseline in Appendix A.5, demonstrating SAC’s superiority for learning adaptive thresholds. The threshold distribution is detailed in

### Appendix A.6.

**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.7.

**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.

### 5.1 Ablation Study

Ablating each of RLShield’s three core modules under the LLaMA-3.1-8B-Instruct backbone (Table 3) reveals that removing dynamic reference selection (-DRS) incurs the steepest drop across all benchmarks, confirming that retrieval-guided rewriting is essential for covering diverse jailbreak

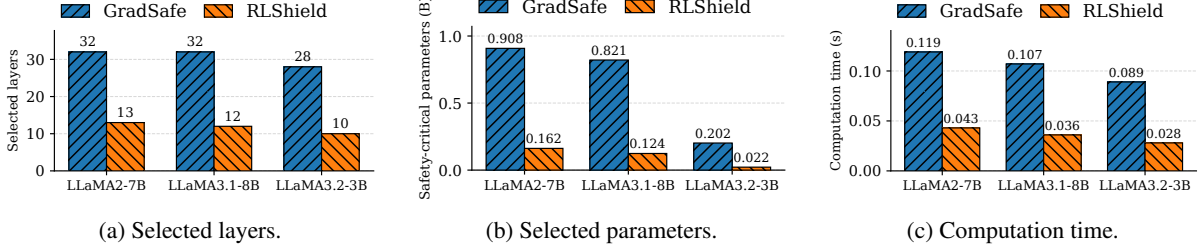


Figure 3: Comparison of selected layers, safety-critical parameters, and computation time across target LLMs.

surface forms; excising critical parameter localization (-CPL) yields the second largest degradation, as unfiltered gradients inject noise into the harmfulness score; and stripping the RL-based adaptive threshold (-AT) further lowers F1, showing that a fixed boundary cannot accommodate prompt-level risk drift. Together, the ablations underscore the synergy of retrieval, localization and reinforcement learning in maintaining robust, agile detection.

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>

## 5.2 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

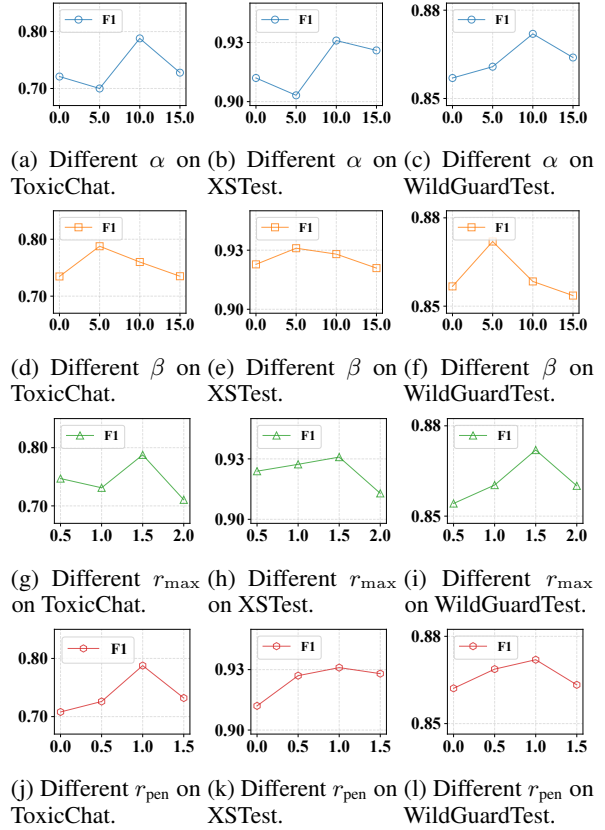


Figure 4: Effect of different  $\alpha$ ,  $\beta$ ,  $r_{\max}$  and  $r_{\text{pen}}$  on three datasets. The y-axis represents F1, and the x-axis indicates different values of  $\alpha$ ,  $\beta$ ,  $r_{\max}$  and  $r_{\text{pen}}$ .

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.

## 562 Limitations

563 Despite the effectiveness of our proposed RLShield,  
564 there are several limitations. First, the safe/unsafe  
565 prompt knowledge base may become outdated as  
566 new jailbreak prompts emerge, making it crucial  
567 to develop efficient methods for its extension and  
568 updating. Second, the evaluation of RLShield has  
569 been conducted exclusively on text-only prompts,  
570 raising questions about its generalizability to multi-  
571 modal inputs, such as vision, audio, and executable  
572 code. Third, our assessment has been limited to  
573 LLaMA- and Qwen-series backbone LLMs, leav-  
574 ing the effectiveness of our approach on other archi-  
575 tectures yet to be investigated. In future work,  
576 we aim to broaden our research to include defenses  
577 against multimodal and executable-code jailbreaks,  
578 as well as apply our methodology to other series of  
579 backbone LLMs.

## 580 References

581 Elisa Bassignana, Valerio Basile, Viviana Patti, and  
582 1 others. 2018. Hurtlex: A multilingual lexicon  
583 of words to hurt. In *CEUR Workshop proceedings*.  
584 CEUR-WS.

585 Guorui Chen, Yifan Xia, Xiaojun Jia, Zhijiang Li, Philip  
586 Torr, and Jindong Gu. 2025. [Llm jailbreak detection  
587 for \(almost\) free!](#) In *Findings of the Association  
588 for Computational Linguistics: EMNLP 2025*, page  
589 5777–5807. Association for Computational Linguistics.  
590

591 Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu  
592 Lian, and Zheng Liu. 2024. Bge m3-embedding:  
593 Multi-lingual, multi-functionality, multi-granularity  
594 text embeddings through self-knowledge distillation.  
595 *ArXiv preprint*.

596 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,  
597 Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,  
598 Akhil Mathur, Alan Schelten, Amy Yang, Angela  
599 Fan, and 1 others. 2024. The llama 3 herd of models.  
600 *arXiv e-prints*.

601 Jessica Maria Echterhoff, Fartash Faghri, Raviteja Vem-  
602 ulapalli, Ting-Yao Hu, Chun-Liang Li, Oncel Tuzel,  
603 and Hadi Pouransari. 2024. Muscle: A model up-  
604 date strategy for compatible llm evolution. In *Find-  
605 ings of the Association for Computational Linguistics:  
606 EMNLP 2024*.

607 Yanlin Feng, Simone Papicchio, and Sajjadur Rahman.  
608 2025. Cypherbench: Towards precise retrieval over  
609 full-scale modern knowledge graphs in the llm era.  
610 In *Proceedings of the 63rd Annual Meeting of the  
611 Association for Computational Linguistics (Volume  
612 1: Long Papers)*.

Peijian Gu, Quan Wang, and Zhendong Mao. 2025. Im-  
prove safety training of large language models with  
safety-critical singular vectors localization. In *Pro-  
ceedings of the 63rd Annual Meeting of the Associa-  
tion 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 Interna-  
tional 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 mod-  
eration tools for safety risks, jailbreaks, and refusals  
of llms. In *Advances in Neural Information Pro-  
cessing 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 Sys-  
tems 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 at-  
tacks 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 ad-  
versarial attacks against aligned language models.  
*ArXiv preprint*.

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 *Find-  
ings of the Association for Computational Linguistics:  
EMNLP 2023*, Singapore.

Michael L Littman. 2015. Reinforcement learning im-  
proves behaviour from evaluative feedback. *Nature*,  
(7553).

Yi Liu, Junzhe Yu, Huijia Sun, Ling Shi, Gelei Deng,  
Yuqi Chen, and Yang Liu. 2024. Efficient detec-  
tion of toxic prompts in large language models. In  
*Proceedings of the 39th IEEE/ACM International  
Conference on Automated Software Engineering*.

669	AI @ Meta Llama Team. 2024. The llama 3 herd of models.	723
670		724
671	Edward Loper and Steven Bird. 2002. <a href="#">NLTK: The natural language toolkit</a> . In <i>Proceedings of the ACL-02 Workshop on Effective Tools and Methodologies for Teaching Natural Language Processing and Computational Linguistics</i> , Philadelphia, Pennsylvania, USA.	725
672		726
673		727
674		728
675		729
676		
677	Sarah Masud, Sahajpreet Singh, Viktor Hangya, Alexander Fraser, and Tanmoy Chakraborty. 2024. Hate personified: Investigating the role of llms in content moderation. <i>ArXiv preprint</i> .	
678		
679		
680		
681	Meta. 2024. meta-llama/llama-3.2-3b-instruct. Hugging Face model card.	
682		
683	Meta. 2025. meta-llama/llama-guard-4-12b. Hugging Face model card.	
684		
685	Microsoft. 2024. Azure ai content safety: Detect and moderate harmful content in text and images. <a href="https://learn.microsoft.com/azure/ai-services/content-safety/">https://learn.microsoft.com/azure/ai-services/content-safety/</a> .	
686		
687		
688		
689	OpenAI. 2024. Moderation api: A tool for content moderation in language models. <a href="https://platform.openai.com/docs/guides/moderation">https://platform.openai.com/docs/guides/moderation</a> .	
690		
691		
692	OpenAI. 2025. Gpt-5 system card. <a href="https://cdn.openai.com/gpt-5-system-card.pdf">https://cdn.openai.com/gpt-5-system-card.pdf</a> .	
693		
694	Sindhu Padakandla, Sadbhavana Babar, Manohar Kaul, and 1 others. 2025. Safequant: Llm safety analysis via quantized gradient inspection. In <i>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)</i> .	
695		
696		
697		
698		
699		
700		
701	Yikang Pan, Liangming Pan, Wenhua Chen, Preslav Nakov, Min-Yen Kan, and William Wang. 2023. <a href="#">On the risk of misinformation pollution with large language models</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , Singapore.	
702		
703		
704		
705		
706		
707	Perspective. 2024. Perspective api: A tool for toxicity detection in online content. <a href="https://perspectiveapi.com/">https://perspectiveapi.com/</a> .	
708		
709		
710	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 <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , Mexico City, Mexico.	
711		
712		
713		
714		
715		
716		
717		
718	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. <i>ArXiv preprint</i> .	
719		
720		
721		
722		
	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 <i>Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2</i> .	723
		724
		725
		726
		727
		728
		729
	Jiong Xiao 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 <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	730
		731
		732
		733
		734
		735
		736
		737
		738
	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 <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> .	739
		740
		741
		742
		743
		744
	Yihan Wang, Zhouxing Shi, Andrew Bai, and Cho-Jui Hsieh. 2024b. Defending llms against jailbreaking attacks via backtranslation. <i>ArXiv preprint</i> .	745
		746
		747
	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 <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	748
		749
		750
		751
		752
		753
		754
		755
		756
	Yueqi Xie, Minghong Fang, Renjie Pi, and Neil Gong. 2024. Gradsafe: Detecting jailbreak prompts for llms via safety-critical gradient analysis. <i>ArXiv preprint</i> .	757
		758
		759
	Zitao Xuan, Xiaofeng Mao, Da Chen, Xin Zhang, Yuhua Dong, and Jun Zhou. 2025. Shieldhead: Decoding-time safeguard for large language models. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> .	760
		761
		762
		763
		764
	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. <i>ArXiv preprint</i> .	765
		766
		767
		768
	Kai Ye, Hongyi Zhou, Jin Zhu, Francesco Quinzan, and Chengchun Shi. 2025. Robust reinforcement learning from human feedback for large language models fine-tuning. <i>ArXiv preprint</i> .	769
		770
		771
		772
	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 <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> .	773
		774
		775
		776
		777
		778

## 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.

### 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 4.

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.

### 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

Table 4: 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

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 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  $\text{top-}k = 3$  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 identified safety-layer set  $\ell_{\text{safe}}$  for each target LLM in Table 5.

Table 5: 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}

#### 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-PCI-E-32GB GPUs. We compute the perplexity and token log-probability mean and variance features

### 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.

857 in the state using LLaMA2-7B-chat-hf. The ac-  
858 tion is a scalar threshold  $a_t \in [T_{\min}, T_{\max}]$  with  
859  $T_{\min} = -0.1$  and  $T_{\max} = 0.5$ , and the decision  
860 rule is  $\hat{y} = 1$  (harmful) if harmfulness  $> a_t$ , other-  
861 wise  $\hat{y} = 0$  (safe). We adopt the standard off-policy  
862 SAC framework with an actor and twin critics, opti-  
863 mized using Adam; the temperature is adjusted via

864 automatic entropy tuning. For the reward follows  
865 Eq.(11), we set the hyperparameters as  $\alpha = 10$ ,  
866  $\beta = 5$ ,  $r_{\max} = 1.5$ , and  $r_{\text{pen}} = 1$ .

#### 867 A.4 Evaluation Metrics

868 Following our main experiments, we report  
869 AUPRC as a threshold-free score-based metric and

### 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.

**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

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

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.

Table 6: F1 (%) comparison between the regression baseline and RLShield across three benchmarks.

Base LLM	Method	ToxicChat	XSTest	WildGuardTest
LLaMA2-7B	Regression	67.2	26.2	77.4
	RLShield	80.9	91.5	87.0
LLaMA3.1-8B	Regression	69.0	30.5	76.5
	RLShield	78.8	93.1	86.8
LLaMA3.2-3B	Regression	62.2	24.1	76.6
	RLShield	71.1	88.3	83.8

## A.5 SAC vs. Regression

Besides SAC, we also implement a linear regression to predict an prompt-specific decision threshold. Concretely, we fit a logistic regression model that combines the base harmfulness score with feature-based calibration:

$$P(\hat{y} = 1 | X, \text{harmfulne}) = \sigma(\text{harmfulness} - X^\top \beta), \quad (17)$$

where  $\sigma(\cdot)$  is the sigmoid function,  $X$  is the concatenation of the same feature fields used in SAC’s state (semantic features, statistical features, and sentiment features), and  $\beta$  contains an intercept term and feature weights. Under the default decision rule

$$\hat{y} = 1 \iff \text{harmfulness} > X^\top \beta, \quad (18)$$

i.e., the predicted threshold equals  $X^\top \beta$ .

We optimize  $\beta$  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. In our implementation, we solve this convex objective with LBFGS (Strong Wolfe line search) on GPU.

We compare this linear regression with our SAC-based adaptive thresholding. As shown in Table 6, the linear baseline is consistently weaker than SAC on all three benchmarks under different backbones, with a particularly clear gap on ToxicChat and

Table 7: 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

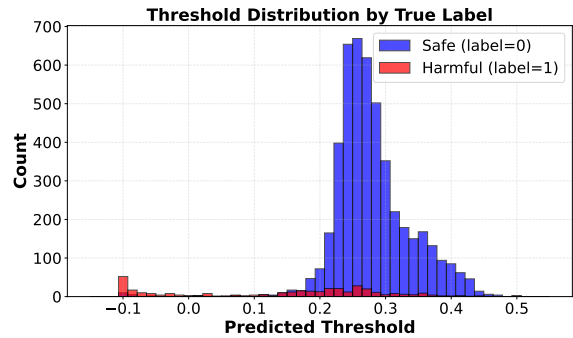
WildGuardTest. This supports the necessity of SAC: a purely regression-based threshold predictor is often too restrictive and cannot capture the complex dependence between the base harmfulness score and the state features, whereas SAC can learn an prompt-specific, complex thresholding policy optimized directly by reward feedback.

### A.6 Threshold Distribution of the SAC Agent

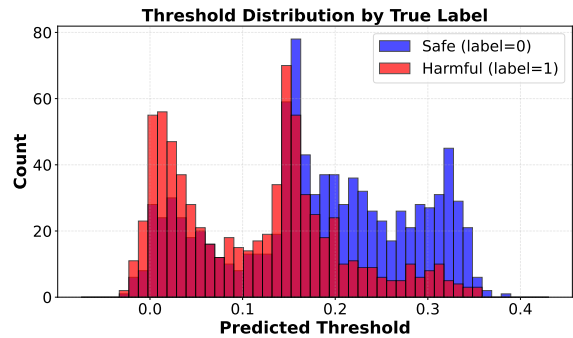
To better understand how the learned policy adapts the decision boundary, we visualize the distribution of the *predicted thresholds* produced by the SAC agent (trained with LLaMA3.1-8B-Instruct as the base model) on three test sets. Fig. 7 plots histograms of the predicted thresholds, stratified by 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.7 Additional Backbones and Model Scales

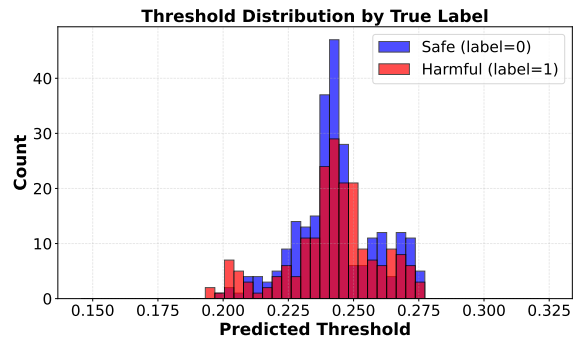
To examine the generality of our method across different model families and parameter scales, we further evaluate it on additional safety-aligned LLM backbones beyond the LLaMA series. Table 7 summarizes the results on three benchmarks. Overall, RLShield consistently yields competitive performance across these diverse backbones, indicating that its effectiveness is not tied to a particular model family or scale.



(a) ToxicChat



(b) WildGuardTest



(c) XSTest

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