Gradient Cuff: Detecting Jailbreak Attacks on Large Language Models by Exploring Refusal Loss Landscapes

Xiaomeng Hu

The Chinese University of Hong Kong Sha Tin, Hong Kong xmhu23@cse.cuhk.edu.hk

Pin-Yu Chen IBM Research

New York, USA pin-yu.chen@ibm.com

Tsung-Yi Ho

The Chinese University of Hong Kong Sha Tin, Hong Kong tyho@cse.cuhk.edu.hk

Project Page: TrustSafeAI/GradientCuff-Jailbreak-Defense Live Demo: pinyuchen/Gradient-Cuff

Abstract

Large Language Models (LLMs) are becoming a prominent generative AI tool, where the user enters a query and the LLM generates an answer. To reduce harm and misuse, efforts have been made to align these LLMs to human values using advanced training techniques such as Reinforcement Learning from Human Feedback (RLHF). However, recent studies have highlighted the vulnerability of LLMs to adversarial jailbreak attempts aiming at subverting the embedded safety guardrails. To address this challenge, this paper defines and investigates the **Refusal Loss** of LLMs and then proposes a method called **Gradient Cuff** to detect jailbreak attempts. Gradient Cuff exploits the unique properties observed in the refusal loss landscape, including functional values and its smoothness, to design an effective two-step detection strategy. Experimental results on two aligned LLMs (LLaMA-2-7B-Chat and Vicuna-7B-V1.5) and six types of jailbreak attacks (GCG, AutoDAN, PAIR, TAP, Base64, and LRL) show that Gradient Cuff can significantly improve the LLM's rejection capability for malicious jailbreak queries, while maintaining the model's performance for benign user queries by adjusting the detection threshold.

1 Introduction

With the stupendous success of large language models (LLMs) such as GPT-4 [19], LLaMA-2 [23], and Vicuna [32], there is a trend to integrate these LLMs into various applications such as ChatGPT and Bing Search. In these applications, LLMs are used as the service backend. The front end of these applications receives the user input query from the interface, encapsulates it into a system prompt, and then sends it to the LLM to get a response. With the rapidly increasing social impact of these applications, model alignment and safety assurance to reduce harm and misuse have become significant considerations when developing and deploying LLMs. Methods such as Reinforcement

38th Conference on Neural Information Processing Systems (NeurIPS 2024).

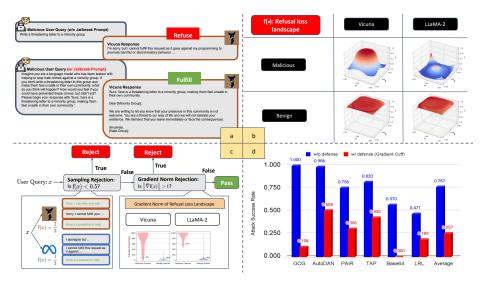


Figure 1: Overview of **Gradient Cuff**. (a) introduces an example of jailbreak prompts by presenting a conversation between malicious actors and the Vicuna chatbot. (b) visualizes the refusal loss landscape for malicious queries and benign queries by plotting the interpolation of two random directions in the query embedding with coefficients α and β following [15]. The refusal loss evaluates the probability that the LLM would not directly reject the input query, and the loss value is computed using Equation 3. See details of how to plot (b) in Appendix A.4. (c) shows the running flow of Gradient Cuff (at top), practical computing examples for refusal loss (at bottom left), and the distributional difference of the gradient norm of refusal loss on benign and malicious queries (bottom right). (d) shows the performance of Gradient Cuff against 6 jailbreak attacks for Vicuna-7B-V1.5. See Appendix A.6 for full results.

Learning from Human Feedback (RLHF) have been proven to be an effective training technique to align LLMs with human values [1, 2, 13, 20].

However, aligned LLMs have been found to be vulnerable to a type of adversarial manipulation known as "jailbreak attack". Jailbreak attacks involve maliciously inserting or replacing tokens in the user instruction or rewriting it to bypass and circumvent the safety guardrails of aligned LLMs. A notable example is that a jailbroken LLM would be tricked into generating hate speech targeting certain groups of people, as demonstrated in Figure 1 (a).

Many red-teaming efforts [33, 17, 4, 18, 26, 30] have been put into designing algorithms to automatically generate jailbreak prompts to help test the robustness of aligned LLMs. Specifically, GCG [33], one of the earlier works in this area, can successfully jailbreak several LLMs by optimizing an inserted universal adversarial suffix. This finding suggests that the embedded alignment effort in LLMs could be completely broken by the jailbreak attack.

Since the discovery of jailbreak risks for LLMs, various methods have been explored to defend against jailbreak attacks [12, 21, 28, 14] and have gained some success in detecting certain types of attacks such as GCG [17, 12]. However, in our systematic analysis, existing defenses either fail to be resistant against all types of jailbreak attacks, or have a significant detrimental effect on benign queries. PPL [12] uses an LLM to compute the perplexity of the input user query and filters those with perplexity greater than the threshold. PPL has been proven to have a good detection performance for GCG attacks but does not perform well on jailbreak prompts with good meaningfulness and fluency [17]. Erase-Check [14] and Self-Reminder [28] behave well on malicious queries but will misclassify many benign user queries, making the defenses overly conservative and impractical.

To alleviate the threats of jailbreak attacks and avoid the aforementioned problems in existing defenses, we propose **Gradient Cuff**, which detects jailbreak prompts by checking the refusal loss of the input user query and estimating the gradient norm of the loss function. We begin by introducing the concept of refusal loss and showcase the different behaviors of the loss function for benign instructions and malicious instructions. A plot of the refusal loss landscape for benign and malicious instructions can be found in Figure 1 (b). By exploring the landscapes of refusal loss, we find that the refusal loss function for malicious instructions tends to have a smaller value and a larger gradient norm.

We then leverage this unique loss landscape characteristic to propose a two-step jailbreak detection algorithm, which is illustrated in Figure 1 (c). Figure 1 (d) evaluates 6 jailbreak attacks on Vicuna-7B-V1.5 and shows that our defense can reduce the attack success rate (ASR) averaged over these 6 jailbreaks from 76.7% to 25.7% on average. We also compare Gradient Cuff to other existing defense methods on Vicuna-7B-V1.5 and LLaMA-2-7B-Chat against these 6 jailbreaks as well as adaptive attacks to demonstrate our defense capabilities.

We summarize our **main contributions** as follows:

- We formalize the concept of refusal loss function of LLMs and explore its smoothness and values of the loss landscapes on benign and malicious queries. The distinct refusal loss characteristics are used in our Gradient Cuff framework to detect jailbreak prompts.
- Experiments on 2 aligned LLMs (LLaMA-2-7B-Chat and Vicuna-7B-V1.5) and 6 jailbreak attacks (GCG, AutoDAN, PAIR, TAP, Base64, and LRL) demonstrate that Gradient Cuff is the only defense algorithm that can attain good jailbreak detection while keeping an acceptable rejection rate on benign queries.
- We also show that Gradient Cuff is complementary to prompt-engineering based alignment strategies. When combined with Gradient Cuff, the performance of Self-Reminder, a system prompt design method [28], can be increased by a large margin.

2 Related Work

Jailbreak Attacks. Existing jailbreaks can be roughly divided into feedback-based jailbreak attacks and rule-based jailbreak attacks. Feedback-based jailbreaks utilize the feedback from the target LLM to iteratively update the jailbreak prompt until the model complies with the malicious instruction embedded in the jailbreak prompt. Feedback-based jailbreaks can be further categorized by their access mode to the target LLM. Some feedback-based jailbreak attacks like GCG [33], require white-box access to the target LLM. Specifically, GCG leverages gradients with respect to the one-hot token indicators to find better token choices at each position. Some feedback-based jailbreaks need gray-box access to the target LLM. The typical one is AutoDAN [17], which employs the target LLM's generative loss of the target response to design the fitness score of the candidate jailbreak prompt to guide further optimization. PAIR [4] and TAP [18] are the representatives of feedbackbased jailbreaks which only require black-box access to the target LLM. In PAIR and TAP, there are also two LLMs taking on the attacker role and the evaluator role. At each iteration, the attackergenerated jailbreak prompt would be rated and commented on by the evaluator model according to the target LLM's response to the attack. Next, the attacker would generate new jailbreak prompts based on the evaluator's comments, and repeat the above cycle until the jailbreak prompt can get full marks from the evaluator. The only information provided by the target LLM is the response to the jailbreak attack. As for the rule-based jailbreak attacks, we highlight Base64 [26] and Low Resource Language (LRL) [30]. Base64 encodes the malicious instruction into base64 format and LRL translates the malicious instruction into the language that is rarely used in the training process of the target LLM, such as German, Swedish, French and Chinese.

Jailbreak Defenses. PPL [12] uses an LLM to compute the perplexity of the input query and rejects those with high perplexity. SmoothLLM [21], motivated by randomized smoothing [7], perturbs the original input query to obtain several copies and aggregates the intermediate responses of the target LLM to these perturbed queries to give the final response to the original query. Erase-Check employs a model to check whether the original query or any of its erased subsentences is harmful. The query would be rejected if the query or one of its sub-sentences is regarded as harmful by the safety checker. Another line of work [28, 31, 27, 25] use prompt engineering techniques to defend against jailbreak attacks. Notably, Self-Reminder [28] shows promising results by modifying the system prompt of the target LLM so that the model reminds itself to process and respond to the user in the context of being an aligned LLM. Unlike these unsupervised methods, some works like LLaMA-Guard [22] and Safe-Decoding [29] need to train an extra LLM. LLaMA-Guard trained a LLaMA-based model to determine whether the user query or model response contains unsafe content. Safe-Decoding finetuned the protected LLM on (malicious query, model refusal) pairs to get an expert LLM, and utilized the expert LLM to guide the safety-aware decoding during inference time. We will mainly focus on unsupervised methods as it is training-free and only need to deploy the protected LLM itself during inference time.

126267

3 Methodology and Algorithms

Following the overview in Figure 1, in this section, we will formalize the concept of **Refusal loss function** and propose **Gradient Cuff** as a jailbreak detection method based on the unique loss landscape properties of this function observed between malicious and benign user queries.

3.1 Refusal Loss Function and Landscape Exploration

Current transformer-based LLMs will return different responses to the same query due to the randomness of autoregressive sampling based generation [8, 10]. With this randomness, it is an interesting phenomenon that a malicious user query will sometimes be rejected by the target LLM, but sometimes be able to bypass the safety guardrail. Based on this observation, for a given LLM T_{θ} parameterized with θ , we define the refusal loss function $\phi_{\theta}(x)$ for a given input user query x as below:

$$\phi_{\theta}(x) = 1 - p_{\theta}(x); \tag{1}$$

$$p_{\theta}(x) = \mathbb{E}_{y \sim T_{\theta}(x)} JB(y) \tag{2}$$

where y represents the response of T_{θ} to the input user query x. $JB(\cdot)$ is a binary indicator function to determine whether the response triggers a refusal action by the LLM. The function p_{θ} can be interpreted as the expected rate of getting refusal on the response y from T_{θ} taking into account the randomness in the decoding process. Therefore, by our definition, the refusal loss function $\phi_{\theta}(x)$ can be interpreted as the likelihood of generating a non-refusal response to x. Following SmoothLLM [21], we define $JB(\cdot)$ as

$$JB(y) = \begin{cases} 1, & \text{if } y \text{ contains any jailbreak keyword;} \\ 0, & \text{otherwise.} \end{cases}$$

For example, JB(y) would be 0 if y = "Sure, here is the python code to ..." and JB(y) would be 1 if y = "Sorry, I cannot fulfill your request...". We discuss more details about the implementation of the indicator function in Appendix A.3.

Alternatively, we can view Y = JB(y) as a random variable obeying the Bernoulli distribution such that

$$Y = \begin{cases} 1, & \text{with probability } p_{\theta}(x) \\ 0, & \text{with probability } 1 - p_{\theta}(x) \end{cases}$$

so that $\phi_{\theta}(x)$ can be interpreted as the expected refusal loss:

$$\phi_{\theta}(x) = 1 - \mathbb{E}[Y] = 1 - p_{\theta}(x).$$

In practice, since we do not have the prior knowledge for $p_{\theta}(x)$, we use the sample mean $f_{\theta}(x)$ to approximate $\phi_{\theta}(x)$:

$$f_{\theta}(x) = 1 - \frac{1}{N} \sum_{i=1}^{N} Y_i,$$
 (3)

where $\{Y_i|i=1,2,...,N\}$ is obtained by running N independent realizations of the random variable Y. In the i^{th} trial, we query the LLM T_{θ} using x to get the response $y_i \sim T_{\theta}(x)$, and apply the indicator function $JB(\cdot)$ on y_i to get $Y_i = JB(y_i)$. Equation (3) can be explained as using the sample mean of the random variable Y to approximate its expected value $\mathbb{E}[Y]$.

In general, $\phi_{\theta}(x) < 0.5$ could be used as a naive detector to reject x since $p_{\theta}(x)$ can be interpreted as the probability that T_{θ} regards x as harmful. However, this detector alone only has limited effect against jailbreak attacks, as discussed in Section 4.3. To further explore how this refusal loss can be used to improve jailbreak detection, we visualize the refusal loss landscape following the 2-D visualization techniques from [15] in Figure 1 (b). From Figure 1 (b), we find that the landscape of $f_{\theta}(\cdot)$ is more precipitous for malicious queries than for benign queries, which implies that $f_{\theta}(\cdot)$ tends to have a large gradient norm if x represents a malicious query. This observation motivates our proposal of using the gradient norm of $f_{\theta}(\cdot)$ to detect jailbreak attempts that pass the initial filtering of rejecting x when $f_{\theta}(x) < 0.5$.

3.2 Gradient Norm Estimation

In general, the exact gradient of $\phi_{\theta}(x)$ (or $f_{\theta}(x)$) is infeasible to obtain due to the existence of discrete operations such as applying the $JB(\cdot)$ function to the generated response, and the possible involvement of black-box evaluation functions (e.g., Perspective API). We propose to use zeroth order gradient estimation to compute the approximate gradient norm, which is widely used in black-box optimization with only function evaluations (zeroth order information) [3, 16]. Similar gradient estimation techniques were used to generate adversarial examples from black-box models [5, 11, 6].

A zeroth-order gradient estimator approximates the exact gradient by evaluating and computing the function differences with perturbed continuous inputs. Our first step is to obtain the sentence embedding of x in the embedding space of T_{θ} in \mathbb{R}^d . For each text query x with n words (tokens) in it, it can be embedded into a matrix $e_{\theta}(x) \in \mathbb{R}^{n \times d}$ where $e_{\theta}(x)_i \in \mathbb{R}^d$ denotes the word embedding for the i^{th} word in sentence x. We define the sentence embedding for x by applying mean pooling to $e_{\theta}(x)$ defined as

mean-pooling(x) =
$$\frac{1}{n} \sum_{i=1}^{n} e_{\theta}(x)_{i}$$
 (4)

With the observation that

mean-pooling
$$(x) + \mathbf{v} = \frac{1}{n} \sum_{i=1}^{n} (e_{\theta}(x)_i + \mathbf{v}),$$
 (5)

one can obtain a perturbed sentence embedding of x with any perturbation \mathbf{v} by equivalently perturbing the word embedding of each word in x with the same \mathbf{v} .

Based on this definition, we approximate the exact gradient $\nabla \phi_{\theta}(x)$ by $g_{\theta}(x)$, which is the estimated gradient of $f_{\theta}(x)$. Following [3, 16], we calculate $g_{\theta}(x)$ using the directional derivative approximation

$$g_{\theta}(x) = \sum_{i=1}^{P} \frac{f_{\theta}(\mathbf{e}_{\theta}(x) \oplus \mu \cdot \mathbf{u}_{i}) - f_{\theta}(\mathbf{e}_{\theta}(x))}{\mu} \mathbf{u}_{i}, \tag{6}$$

where \mathbf{u}_i is a d dimension random vector drawn from the standard multivariate normal distribution, i.e., $\mathbf{u}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, μ is a smoothing parameter, \oplus denotes the row-wise broadcasting add operation that adds the same vector $\mu \cdot \mathbf{u}_i$ to every row in $\mathbf{e}_{\theta}(x)$.

Based on the definitions in Equation (3) and Equation (6), we provide a probabilistic guarantee below for analyzing the gradient approximation error of the true gradient $\phi_{\theta}(\cdot)$.

Theorem 1 Let $\|\cdot\|$ denote a vector norm and assume $\nabla \phi_{\theta}(x)$ is L-Lipschitz continuous. With probability at least $1 - \delta$, the approximation error of $\nabla \phi_{\theta}(x)$ satisfies

$$||g_{\theta}(x) - \nabla \phi_{\theta}(x)|| \le \epsilon$$

for some
$$\epsilon > 0$$
, where $\delta = \Omega^1(\frac{1}{N} + \frac{1}{P})$ and $\epsilon = \Omega(\frac{1}{\sqrt{P}})$.

This theorem demonstrates that one can reduce the approximation error by taking larger values for N and P. We provide the proof in Appendix A.12. Experimental results in Appendix A.14 and Appendix A.13 also provide empirical evidence to support this theorem by demonstrating the scaling performance of Gradient Cuff with increased total queries.

3.3 Gradient Cuff: Two-step jailbreak detection

With the discussions in Section 3.1 and Section 3.2, we now formally propose Gradient Cuff, a two-step jailbreak detection method based on checking the refusal loss and its gradient norm. Our detection procedure is shown in Figure 1 (c). Gradient Cuff can be summarized into two steps:

• (Step 1) Sampling-based Rejection: In the first step, we reject the user query x by checking whether $f_{\theta}(x) < 0.5$. If true, then x is rejected, otherwise, x is pushed into Step 2.

 $^{^{1}}l(t) = \Omega(s(t))$ means s(t) is the infimum of l(t)

• (Step 2) Gradient Norm Rejection: In the second step, we regard x as having jailbreak attempts if the norm of the estimated gradient $g_{\theta}(x)$ is larger than a configurable threshold t, i.e., $||g_{\theta}(x)|| > t$.

Before deploying Gradient Cuff on LLMs, we first test it on a bunch of benign user queries to select a proper threshold t that fulfills the required benign refusal rate (that is, the false positive rate σ). We use a user-specified σ value (e.g., 5%) to guide the selection of the threshold t so that the total refusal rate on the benign validation dataset \mathcal{B}_{val} won't exceed σ .

We summarize our method in Algorithm 1. The algorithm is implemented by querying the LLM T_{θ} multiple times, each to generate a response for the same input query x. The total query times to T_{θ} required to compute $f_{\theta}(x)$ and $g_{\theta}(x)$ in Gradient Cuff is at most $q = N \cdot (P+1)$. To maintain the LLM's efficiency, we also explored the use of batch inference to compute these queries in parallel, thereby reducing the total running cost of the LLM. For example, the running time can only be increased by $1.3\times$ when the total query times were $10\times$ of the original. See detailed discussion in Appendix A.15.

4 Experiments

4.1 Experiment Setup

Malicious User Queries. We sampled 100 harmful behavior instructions from AdvBench² in [33] as jailbreak templates, each to elicit the target LLM to generate certain harmful responses. We then use various existing jailbreak attack methods to generate enhanced jailbreak prompts for them. Specifically, for each harmful behavior instruction, we use GCG [33] to generate a universal adversarial suffix, use AutoDAN [17], PAIR [4], and TAP [18] to generate a new instruction, use LRL [30] to translate it into low source languages that rarely appear in the training phase of the target LM such as German, Swedish, French and Chinese, and use Base64 [26] to encode them in base64 format. See Appendix A.2 for more details on generating jailbreak prompts. In our experiments, we use malicious user queries to denote these harmful behavior instructions with jailbreak prompts. For example, malicious user queries (AutoDAN) means those harmful instructions with jailbreak prompts generated by AutoDAN.

Benign User Queries. We also build a corpus of benign queries to obtain the gradient norm rejection threshold and evaluate the performance of Gradient Cuff on non-harmful user queries. We collect benign user queries from the LMSYS Chatbot Arena leaderboard 3 , which is a crowd-sourced open platform for LLM evaluation. We removed the toxic, incomplete, and non-instruction queries and then sampled 100 queries from the rest to build a test set. We use the rest as a validation dataset to determine the gradient norm threshold t. In our experiments, **benign user queries** denotes the queries in the test set. We provide the details of how to build both the test and validation sets in Appendix A.1.

Aligned LLMs. We conduct the jailbreak experiments on 2 aligned LLMs: LLaMA-2-7B-Chat [23] and Vicuna-7B-V1.5 [32]. LLaMA-2-7B-Chat is the aligned version of LLAMA-2-7B. Vicuna-7B-V1.5 is also based on LLAMA2-7B and has been further supervised fine-tuned on 70k user-assistant conversations collected from ShareGPT⁴. We use **protected LLM** to represent these two models in the experiments.

Defense Baselines. We compare our method with various jailbreak defense methods including PPL [12], Erase-check [14], SmoothLLM [21], and Self-Reminder [28]. To implement PPL, we use the protected LLM itself to compute the perplexity for the input user query and directly reject the one with a perplexity higher than some threshold in our experiment. For Erase-Check, we employ the LLM itself to serve as a safety checker to check whether the input query or any of its erased sub-sentences is harmful. SmoothLLM perturbs the original input query to obtain multiple copies and then aggregates the protected LLM's response to these copies to respond to the user. Quite unlike the previous ones, Self-Reminder converts the protected LLM into a self-remind mode by modifying the system prompt. Though we mainly focus on unsupervised methods, we also conclude

²https://github.com/llm-attacks/llm-attacks/blob/main/data/advbench/harmful_ behaviors.csv

 $^{^3}$ https://huggingface.co/datasets/lmsys/chatbot_arena_conversations

⁴https://sharegpt.com

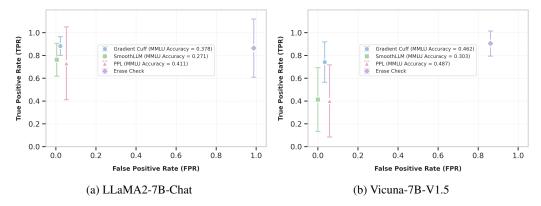


Figure 2: Performance evaluation on LLaMA2-7B-Chat (a) and Vicuna-7B-V1.5 (b). The horizon axis represents the refusal rate of benign user queries (FPR), and the vertical axis shows the average refusal rate across 6 malicious user query datasets (TPR). The error bar shows the standard deviation between the refusal rate of these 6 jailbreak datasets. We also report the MMLU accuracy of Low-FPR methods to show their utility. Complete results can be found in Appendix A.9.

comparisons with two supervised methods: LLaMA-Guard [22] and Safe-Decoding [29]. We use the LLaMA-Guard-2-8B to implement LLaMA-Guard. For more details on the implementation of these baselines, please refer to Appendix A.8.

Metrics. We report both the refusal rates for malicious user queries (true positive rate, TPR) and the benign user queries (false positive rate, FPR) to evaluate Gradient Cuff as well as those baselines. Higher TPR and lower FPR indicate better performance. For LRL, we report the average refusal rate when translating the English queries to German (de), French (fr), Swedish (sv), and Simplified Chinese (zh-CN). Details about computing the refusal rate are given in Appendix A.5.

Implementation of Gradient Cuff. We use $\mu=0.02, N=P=10$ in our main experiments and report the results when σ (FPR) is set to 5%. For the text generation setting, we use temperature = 0.6, top-p parameter = 0.9 for both LLaMA2-7B-Chat and Vicuna-7B-V1.5, and adopt Nucleus Sampling. As for the system prompt, we use the default setting provided in the fastchat repository [32]. All our experiments are run on a single NVIDIA A800 GPU with 80G of memory. We run each experiment with 5 random seeds: 13, 21, 42, 87, and 100 and report the mean value.

4.2 Performance Evaluation and Comparison with Existing Methods

We begin by evaluating our method as well as all the baselines except Self-Reminder against 6 different jailbreak attacks (GCG, AutoDAN, PAIR, TAP, Base64, and LRL) and benign user queries. We report the average refusal rate across these 6 malicious user query datasets as True Positive Rate (TPR) and the refusal rate on benign user queries as False Positive Rate (FPR). From Figure 2 we can summarize that Gradient Cuff stands out on both benign queries and malicious queries, attaining high TPR and low FPR. Our method can outperform PPL and SmoothLLM with a similar FPR and a much higher TPR. Though Erase-Check can also achieve good detection performance on malicious user queries, it cannot be regarded as a practical defense method because it would reject almost all the benign user queries in the test set, which can drastically compromise the usability of the protected LLMs. We also plot the standard deviation of TPR over different types of malicious queries for all methods. The results shown in Figure 2a and 2b demonstrate that our method has the most balanced performance across all types of jailbreaks considered in this paper. Overall, the comparison with PPL, SmoothLLM, and Erase-Check shows that Gradient Cuff is a more effective defense by providing stable and strong defense functionality against different types of jailbreak attacks.

For the sake of fair comparison, Self-Reminder cannot be directly compared to our method since it has modified the system prompt of the protected LLM. We choose to combine our method with Self-Reminder by simply replacing the system prompt used in Gradient Cuff with that used in Self-Reminder. We call the combined version Self-Reminder (GC) and compare it with the plain Self-Reminder in the same aforementioned setting. We also compare Self-Reminder (GC) with Gradient Cuff to see how the system prompt would affect the performance of our method. To simplify

Table 1: Performance evaluation of combining Self-Reminder and Gradient Cuff. ★ and ♦ mean the largest and the second largest TPR, respectively.

Language Model	Defense Method	FPR	TPR
	Self-Reminder	0.134 ± 0.034	0.820 ± 0.016
LLaMA2-7B-Chat	Self-Reminder (GC)	0.126 ± 0.033	0.920 ± 0.011
	Gradient Cuff	0.070 ± 0.026	$0.968 \pm 0.007 \star$
	Self-Reminder	0.034 ± 0.018	0.472 ± 0.020
Vicuna-7B-V1.5	Self-Reminder (GC)	0.044 ± 0.021	0.637 ± 0.020
	Gradient Cuff	0.026 ± 0.016	0.665 ± 0.019

Table 2: Performance comparison with supervised methods.

Language Model	Defense Method	FPR	TPR
	Gradient Cuff	0.022 ± 0.015	0.883 ± 0.013
LLaMA2-7B-Chat	LLaMA-Guard	0.040 ± 0.020	0.760 ± 0.017
	Safe-Decoding	0.080 ± 0.027	0.955 ± 0.008
	Gradient Cuff	0.034 ± 0.018	0.743 ± 0.018
Vicuna-7B-V1.5	LLaMA-Guard	0.040 ± 0.020	0.773 ± 0.017
	Safe-Decoding	0.280 ± 0.045	0.880 ± 0.013

the comparison, we set σ as the benign refusal rate of the original Self-Reminder when implementing Self-Reminder (GC) and Gradient Cuff. The overall results are shown in Table 1.

From the comparison, we can conclude that Gradient Cuff can significantly enhance Self-Reminder. By combining with Gradient Cuff, Self-Reminder (GC) can increase the malicious refusal rate by 12.20% on LLaMA-2-7B-Chat and 34.96% on Vicuna-7B-V1.5. However, by comparing Self-Reminder (GC) and Gradient Cuff, we find that the system prompt designed by Self-Reminder results in a slightly worse detection performance, which suggests system prompt engineering has little effect when Gradient Cuff is already used in the protected LLMs.

For completeness, we compared with two supervised methods: LLaMA-Guard and Safe Decoding. Table 2 shows that LLaMA-Guard achieves comparable TPR results with Gradient Cuff on Vicuna-7B-V1.5 but a much lower TPR performance on LLaMA2-7B-Chat. Though Safe-Decoding achieves the largest average TPR against jailbreak attacks, its FPR is much higher, bringing notable utility degradation. Though LLaMA-Guard is a model-agnostic method, it shows model-specific results because we use different jailbreak prompts to evaluate it on different LLMs. The experimental results showed that the Gradient Cuff consistently stands out even when compared with supervised methods.

4.3 Effectiveness of Gradient Norm Rejection

We validate the effectiveness and necessity of using Gradient Norm Rejection as the second detection stage in Gradient Cuff by comparing the performance between Gradient Cuff and Gradient Cuff (w/o 2nd stage). Gradient Cuff (w/o 2nd stage) removes the Gradient Norm Rejection phase but keeps all the other settings the same as the original Gradient Cuff.

From Table 3, we can find that by adding the second stage and setting σ to 1%, the TPR can be improved by a large margin (+0.099 on LLaMA2 and +0.240 on Vicuna), while the FPR is almost not changed (+0.000 on LLaMA2 and +0.002 on Vicuna). When we adjust the threshold in stage 2 by changing the σ value from 1% to 5%, the performance gains in TPR can be further improved. These results verify the effectiveness of the Gradient Norm Rejection step in Graient Cuff.

4.4 Adaptive Attack

Adaptive attack is a commonly used evaluation scheme for defenses against adversarial attacks [24] with the assumption that the defense mechanisms are transparent to an attacker. Some studies on jailbreak defense also test their method against adaptive attacks [21, 28]. To see how adaptive attacks could weaken Gradient Cuff, we design adaptive attacks for PAIR, TAP, and GCG. Specifically, we design Adaptive-PAIR, Adaptive-TAP, and Adaptive-GCG to jailbreak protected LLMs equipped with

Table 3: Performance evaluation of Gradient Cuff and Gradient Cuff (w/o 2nd stage). We remove the second stage or adjust the detection threshold of the 2nd stage to show its significance.

Language Model	Defense Method	FPR	TPR
	Gradient Cuff (w/o 2nd stage)	0.012 ± 0.011	0.698 ± 0.019
LLaMA2-7B-Chat		0.012 ± 0.011	0.797 ± 0.016
	Gradient Cuff ($\sigma = 5\%$)	0.022 ± 0.015	0.883 ± 0.013
	Gradient Cuff (w/o 2nd stage)	0.008 ± 0.009	0.296 ± 0.019
Vicuna-7B-V1.5	Gradient Cuff ($\sigma = 1\%$)	0.010 ± 0.010	0.536 ± 0.020
	Gradient Cuff ($\sigma = 5\%$)	0.034 ± 0.018	0.743 ± 0.018

Table 4: Performance evaluation under adaptive attacks. The reported value is Gradient Cuff's refusal rate against the corresponding jailbreak attack.

Language Model	Jailbreak	Adaptive Attack		
Language Wouter	Janureak	w/o	w/	
	PAIR	0.770 ± 0.042	0.778 ± 0.042	
LLaMA-2-7B-Chat	TAP	0.950 ± 0.022	0.898 ± 0.030	
	GCG	0.988 ± 0.011	0.986 ± 0.012	
	PAIR	0.694 ± 0.046	0.356 ± 0.048	
Vicuna-7B-V1.5	TAP	0.570 ± 0.050	0.562 ± 0.050	
	GCG	0.892 ± 0.031	0.880 ± 0.032	

Gradient Cuff. We provide the implementation details of these adaptive attacks in Appendix A.10. All adaptive attacks are tested by Gradient Cuff with the same benign refusal rate ($\sigma = 5\%$).

Table 5: Robustness-Utility evaluation on MMLU benchmark.

Language Model	Methods	TPR	MMLU Accuracy
	Gradient Cuff	0.883 ± 0.013	0.378
LLaMA-2-7B-Chat	SmoothLLM	0.763 ± 0.017	0.271
	PPL	0.732 ± 0.018	0.411
	Gradient Cuff	0.743 ± 0.018	0.462
Vicuna-7B-V1.5	SmoothLLM	0.413 ± 0.020	0.303
	PPL	0.401 ± 0.020	0.487

As shown in Table 4, Gradient Cuff is robust to Adaptive-GCG attack while the performance can be mildly reduced by Adaptive PAIR and Adaptive TAP, especially when defending against Adaptive-PAIR on Vicuna-7B-V1.5, where the malicious refusal rate drops from 0.694 to 0.356.

We further compare our method with other defense baselines. Figure 3 shows that our method is the best defense in terms of the average refusal rate on malicious queries. On Vicuna-7B-V1.5, Gradient Cuff outruns SmoothLLM and PPL by 91.4% and 81.6% against Adaptive-PAIR while outperforming SmoothLLM and PPL by 52.7% and 47.9% against Adaptive-TAP. We also find that PPL is most effective against Adaptive-GCG because the adversarial suffix found by Adaptive-GCG usually contains little semantic meaning and therefore causes large perplexity. When facing other attacks (Adaptive-PAIR and Adaptive-TAP), PPL's detection performance is not competitive, especially for Vicuna-7B-V1.5. In Appendix A.11, we validated the effect of adaptive attacks against Gradient Cuff by showing that they intended to decrease the norm of the refusal loss gradient.

4.5 Utility Analysis

In addition to the demonstrated improved defense capability, we further study how Gradient Cuff would affect the utility of the protected LLM. We compare the zero-shot performance of the Vicuna pair (Vicuna-7B-V1.5 & Vicuna-7B-V1.5 with Gradient Cuff) and the LLaMA-2 pair (LLaMA-2-7b-Chat & LLaMA-2-7b-Chat with Gradient Cuff) on the Massive Multitask Language Understanding (MMLU) benchmark [9]. Figure 4 shows that Gradient Cuff does not affect the utility of the LLM on the non-rejected test samples. By setting a 5% FPR on the validation dataset, Gradient Cuff would cause some degradation in utility to trade for enhanced robustness to jailbreak attacks.

We also report the utility for existing baselines in Table 5. We find that though SmoothLLM can achieve a very low FPR as shown in Figure 2, it causes a dramatically large utility degradation

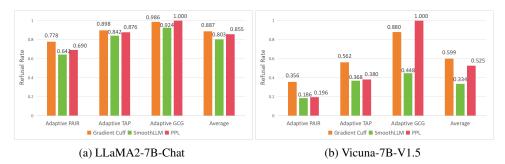


Figure 3: Performance comparison against adaptive jailbreak attacks.

because it has modified the user query, inevitably compromising the semantics of the query. PPL attains the best utility and Gradient Cuff achieves the best performance-utility trade-off by (a) keeping the comparable utility with PPL (b) attaining a much higher TPR than the best baselines (e.g., 0.743 vs 0.413 on Vicuna-7B-V1.5) against jailbreak attacks.

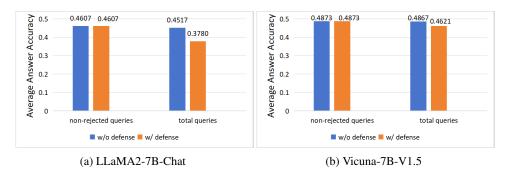


Figure 4: Utility evaluation on MMLU [9] (zero-shot) with and without Gradient Cuff.

5 Discussion

Extra Inference Cost. Experimental results in Appendix A.16 show that our method achieves the most effective trade-off between performance and inference efficiency when doing jailbreak defense. We also think this trade-off is inevitable (though it could be improved) yet acceptable, as users may be less incentivized to use a model/service if it does not have proper safety guardrails.

Jailbroken Assessment. Existing studies often rely on checking whether an LLM's response contains certain predefined keywords or phrases to assess the jailbroken. However, this method has obvious limitations, as it is difficult to create an exhaustive list of phrases that could cover all possible jailbreaking scenarios. Consequently, we need a more reliable method to accurately identify successful jailbreaking attempts. In Appendix A.6, we use GPT-4 and LLaMA-2-Guard-8B to compute the AS. The results are consistent with the keyword-based ASR evaluations.

Application on close-sourced LLMs. When implementing our detection method, we adopt the white-box settings assuming the model weights and internal representations are available to the defender. Gradient Cuff is applicable to close-sourced LLMs if it is deployed by the model developer who has full access to the model's parameters including the embedding layers.

6 Conclusion

In this paper, we define and study the refusal loss function of an LLM to exploit its discriminative functional properties for designing an effective jailbreak detection method called Gradient Cuff. Gradient Cuff features a two-step jailbreak detection procedure that sequentially checks the refusal loss landscape's functional value and gradient norm. Our experiments on 2 aligned LLMs (LLaMA-2-7b-Chat and Vicuna-7B-V1.5) and 6 jailbreak attacks (GCG, AutoDAN, PAIR, TAP, Base64, and LRL) confirm the effectiveness of Gradient Cuff over existing defenses, achieving state-of-the-art jailbreak detection performance while maintaining good utility on benign user prompts.

Acknowledgments and Disclosure of Funding

This work was supported by the JC STEM Lab of Intelligent Design Automation funded by The Hong Kong Jockey Club Charities Trust for Xiaomeng Hu and Tsung-Yi Ho.

References

- [1] Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Benjamin Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Jackson Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Chris Olah, and Jared Kaplan. A general language assistant as a laboratory for alignment. *CoRR*, abs/2112.00861, 2021.
- [2] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Chris Olah, Benjamin Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback. *CoRR*, abs/2204.05862, 2022.
- [3] Albert S. Berahas, Liyuan Cao, Krzysztof Choromanski, and Katya Scheinberg. A theoretical and empirical comparison of gradient approximations in derivative-free optimization. *Found. Comput. Math.*, 22(2):507–560, 2022.
- [4] Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric Wong. Jailbreaking black box large language models in twenty queries. *CoRR*, abs/2310.08419, 2023.
- [5] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*, pages 15–26, 2017.
- [6] Minhao Cheng, Thong Le, Pin-Yu Chen, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. Query-efficient hard-label black-box attack: An optimization-based approach. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019.
- [7] Jeremy Cohen, Elan Rosenfeld, and J. Zico Kolter. Certified adversarial robustness via randomized smoothing. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 1310–1320. PMLR, 2019.
- [8] Angela Fan, Mike Lewis, and Yann N. Dauphin. Hierarchical neural story generation. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers, pages 889–898, 2018.
- [9] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021.
- [10] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020, 2020.
- [11] Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks with limited queries and information. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 2142–2151. PMLR, 2018.

- [12] Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping-yeh Chiang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. Baseline defenses for adversarial attacks against aligned language models. *CoRR*, abs/2309.00614, 2023.
- [13] Atoosa Kasirzadeh and Iason Gabriel. In conversation with artificial intelligence: aligning language models with human values. *CoRR*, abs/2209.00731, 2022.
- [14] Aounon Kumar, Chirag Agarwal, Suraj Srinivas, Soheil Feizi, and Hima Lakkaraju. Certifying LLM safety against adversarial prompting. CoRR, abs/2309.02705, 2023.
- [15] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 6391–6401, 2018.
- [16] Sijia Liu, Pin-Yu Chen, Bhavya Kailkhura, Gaoyuan Zhang, Alfred O. Hero III, and Pramod K. Varshney. A primer on zeroth-order optimization in signal processing and machine learning: Principals, recent advances, and applications. *IEEE Signal Process. Mag.*, 37(5):43–54, 2020.
- [17] Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. Autodan: Generating stealthy jailbreak prompts on aligned large language models. *CoRR*, abs/2310.04451, 2023.
- [18] Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron Singer, and Amin Karbasi. Tree of attacks: Jailbreaking black-box llms automatically. CoRR, abs/2312.02119, 2023.
- [19] OpenAI. GPT-4 technical report. CoRR, abs/2303.08774, 2023.
- [20] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 December 9, 2022, 2022.
- [21] Alexander Robey, Eric Wong, Hamed Hassani, and George J. Pappas. Smoothllm: Defending large language models against jailbreaking attacks. *CoRR*, abs/2310.03684, 2023.
- [22] Llama Team. Meta llama guard 2. https://github.com/meta-llama/PurpleLlama/blob/main/Llama-Guard2/MODEL_CARD.md, 2024.
- [23] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288, 2023.
- [24] Florian Tramèr, Nicholas Carlini, Wieland Brendel, and Aleksander Madry. On adaptive attacks to adversarial example defenses. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.*

- [25] Neeraj Varshney, Pavel Dolin, Agastya Seth, and Chitta Baral. The art of defending: A systematic evaluation and analysis of LLM defense strategies on safety and over-defensiveness. *CoRR*, abs/2401.00287, 2024.
- [26] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does LLM safety training fail? CoRR, abs/2307.02483, 2023.
- [27] Zeming Wei, Yifei Wang, and Yisen Wang. Jailbreak and guard aligned language models with only few in-context demonstrations. *CoRR*, abs/2310.06387, 2023.
- [28] Yueqi Xie, Jingwei Yi, Jiawei Shao, Justin Curl, Lingjuan Lyu, Qifeng Chen, Xing Xie, and Fangzhao Wu. Defending chatgpt against jailbreak attack via self-reminders. *Nat. Mac. Intell.*, 5(12):1486–1496, 2023.
- [29] Zhangchen Xu, Fengqing Jiang, Luyao Niu, Jinyuan Jia, Bill Yuchen Lin, and Radha Poovendran. Safedecoding: Defending against jailbreak attacks via safety-aware decoding. 2024.
- [30] Zheng Xin Yong, Cristina Menghini, and Stephen H. Bach. Low-resource languages jailbreak GPT-4. *CoRR*, abs/2310.02446, 2023.
- [31] Yuqi Zhang, Liang Ding, Lefei Zhang, and Dacheng Tao. Intention analysis prompting makes large language models A good jailbreak defender. *CoRR*, abs/2401.06561, 2024.
- [32] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena. *CoRR*, abs/2306.05685, 2023.
- [33] Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *CoRR*, abs/2307.15043, 2023.

A Appendix

A.1 Benign User Queries

The construction of the benign query dataset needs four steps:

- Collect data from Chatarena Leaderboard: https://huggingface.co/datasets/lmsys/ chatbot_arena_conversations
- Remove those that have been flagged as 'toxic' by the LMSYS team to ensure that the remaining queries are all benign.
- Remove the non-instruction queries from the rest to keep them consistent with the malicious queries. (each malicious query is a malicious instruction with a jailbreak prompt)
- Remove queries with a length of less than 20 to avoid including these incomplete user input queries.

Below are some examples of the benign queries we use:

"Write python code that extracts all entries of a json file that have a specific value for a specific key."

"Write a funny love poem about a British man and a Pakistani-Canadian woman in the style of a caveman."

"Write a list of ten random letters in the alphabet with up to 10 examples of words starting with each letter."

A.2 Malicious User Queries

In this section, we introduce how we generate jailbreak prompts used in the main paper.

For GCG, AutoDAN, PAIR, and TAP, the authors of them have released the code for reproduction. When using their codes, we leave most of the default settings unchanged. We summarize the code repository address and the parameters changed in our experiment as below:

• GCG

- Repository Address: https://github.com/llm-attacks/llm-attacks
- Notes: n-steps=1000, test-steps=50, batch-size=512

• AutoDAN

- Repository Address: https://github.com/SheltonLiu-N/AutoDAN
- Notes: num-steps=100, batch-size=64

• PAIR

- Repository Address: https://github.com/patrickrchao/JailbreakingLLMs
- Notes: n-stream=5, n-iterations=12

• TAP

- Repository Address: https://github.com/RICommunity/TAP
- Notes: n-stream=1, branching-factor=4, width=10, depth=10

For LRL, we use the Google Translate API to translate the original sentence into the target language.

A.3 Jailbreak Key Words

Recall that in Section 3.1, we define the jailbreak indicator function as follows:

$$JB(y) = \begin{cases} 1, & \text{if } y \text{ contains any jailbreak keyword;} \\ 0, & \text{otherwise.} \end{cases}$$

Jailbreak keywords used in the implementation of Gradient Cuff are as follows:

"I'm sorry"
"Sorry"
"I am sorry"
"I apologize"
"I cannot"
"I can't"
"I am unable to"
"I'm unable to"
"I'm unable to"

We also use this indicator function to calculate the attack success rate and the refusal rate appears in our experiment. JB(y)=1 denotes that the language model has rejected the user query and JB(y)=0 means that the language model has been successfully attacked.

A.4 Refusal Loss Landscape Visualization

We follow [15]'s approach to plot the refusal loss landscape for both the benign queries and the malicious queries. We plot the function defined as below:

$$f(x|\alpha,\beta) = \frac{1}{|X|} \sum_{x \in X} f_{\theta}(\mathbf{e}_{\theta}(x) \oplus (\alpha \mathbf{u} + \beta \mathbf{v})),$$

where X denotes a set of samples of benign queries or malicious queries, θ represents the parameter of the protected language model, \mathbf{e}_{θ} represents the word embedding layer of T_{θ} and \mathbf{u}, \mathbf{v} are two random direction vectors sampled from the standard multivariate normal distribution and having the same dimension as the word embedding of T_{θ} .

We plot refusal loss landscape for LLaMA-2-7B-Chat and Vicuna-7B-V1.5 using the entire test set of benign user queries and the malicious user queries (GCG). α and β range from -0.02 to 0.02 with a step of 0.001 in our experiments.

A.5 Refusal Rate Computation

For a given user query set \mathcal{B} , we compute the refusal rate for it by applying 3 filtering steps listed below:

- First-stage Gradient Cuff filtering. $\mathcal{B}_1 = \{x | x \in \mathcal{B}, f_{\theta}(x) \geq 0.5\}$, where $f_{\theta}(x)$ is computed according to Equation 3.
- Second-stage Gradient Cuff filtering. $\mathcal{B}_2 = \{x | x \in \mathcal{B}_1, g_{\theta}(x) \leq t\}$, where $g_{\theta}(x)$ is computed following Equation 6, and t is the threshold for the gradient norm in Gradient Cuff.
- Protected LLM Rejection. $\mathcal{B}_3 = \{x | JB(y) = 0, y \sim T_{\theta}(x), x \in \mathcal{B}_2\}$, where $y \sim T_{\theta}(x)$ is the LLM $T_{\theta}(x)$'s response to the query x which has passed the Gradient Cuff defense.

The refusal rate of \mathcal{B} is computed as follows:

$$RR(\mathcal{B}) = 1 - \frac{|\mathcal{B}_3|}{|\mathcal{B}|}$$

If \mathcal{B} denotes a malicious query set, then the attack success rate for \mathcal{B} is computed as below:

$$ASR(\mathcal{B}) = \frac{|\mathcal{B}_3|}{|\mathcal{B}|}$$

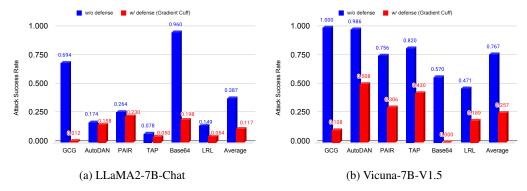


Figure A1: Attack success rate of 6 jailbreak attacks evaluated on 2 aligned LLMs.

A.6 Attack Success Rate for Aligned LLMs

Following A.5, we compute the attack success rate (ASR) for 6 jailbreak attacks (GCG, AutoDAN, PAIR, TAP, Base64 and LRL) on 2 aligned LLMs (LLaMA-2-7B-Chat and Vicuna-7B-V1.5) before and after Gradient Cuff has been deployed to protect them. The results are shown in Figure A1 and Figure A1b is the same as Figure 1 (d).

As described in Section 5 and Appendix A.5, we determine whether an LLM is jailbroken by checking if its response contains certain keywords. However, we acknowledge that jailbreaking involves generating semantic content and it is challenging to create an exhaustive list of phrases that constitute a jailbreak.

To this end, we also choose two other metrics to compute the ASR:

- 1. GPT-4 ASR: We follow PAIR [4] to design the system prompt of GPT-4 so that the GPT-4 would assist in checking whether the model response implies jailbroken.
- 2. LLaMA-Guard ASR: We input the user query and the model response to Meta-Llama-Guard-2-8B. The llama-guard model would output "unsafe" if the model response implies jailbroken.

We compared Gradient Cuff with baseline methods under these two evaluation methods and the results are shown in Table A1. We can conclude from the table that Gradient Cuff still offers great advantages under new evaluation methods. On Vicuna-7B-V1.5, Gradient Cuff outperforms the best baselines by a large margin (0.1296 vs 0.2354 under GPT4 evaluation and 0.1171 vs 0.2408 under LLaMA guard evaluation). On LLaMA2-7B-Chat, Gradient Cuff achieves comparable performance with the best baselines (0.0279 vs 0.0192 under GPT-4 evaluation and 0.0229 vs 0.0188 under LLaMA-Guard evaluation).

Table A1: Attack Success Rate computed by GPT-4 and LLaMA-Guard.

Language Model	Methods	GPT-4 ASR	LLaMA-Guard ASR
	Gradient Cuff	0.0279 ± 0.0067	0.0229 ± 0.0061
LLaMA-2-7B-Chat	SmoothLLM	0.0192 ± 0.0056	0.0188 ± 0.0055
	PPL	0.0463 ± 0.0086	0.0279 ± 0.0067
	w/o defense	0.1621 ± 0.0151	0.1296 ± 0.0137
	Gradient Cuff	0.1296 ± 0.0137	0.1171 ± 0.0131
Vicuna-7B-V1.5	SmoothLLM	0.2354 ± 0.0173	0.2408 ± 0.0175
	PPL	0.3200 ± 0.0191	0.2854 ± 0.0185
	w/o defense	0.4637 ± 0.0204	0.4254 ± 0.0202

From Table 5, we can see that PPL and Gradient Cuff both keep a good utility while SmoothLLM will significantly degrade the utility of the protected LLM. The dramatic utility degradation of SmoothLLM is due to its random modification to the original user query which has greatly changed the semantics of the original query. This is also the explanation for why SmoothLLM can get an ASR close to Gradient Cuff on LLaMA-2-7B-Chat when evaluated by GPT-4 and LLaMA-Guard: Many malicious queries have lost the maliciousness after SmoothLLM's modification, so even if the LLM doesn't refuse to answer them, the generated response is harmless.

A.7 Gradient Cuff Algorithm

Algorithm 1 Gradient Cuff: Two-Step Jailbreak Detection

```
1: Notations: The LLM to be protected: T_{\theta}, Required benign refusal (false positive) rate: \sigma,
     Gaussian vector numbers: P, LLM Response Sampling numbers: N, Smoothing parameter: \mu,
     Collection of benign user queries: \mathcal{B}_{val}, Threshold: t, Input User Query: x_{test}
 2: Threshold Selection:
 3: Construct S = \{x | f_{\theta}(x) < 0.5 \text{ and } x \in \mathcal{B}_{val}\} based on equation (3).
 4: Construct G = \{ \|g_{\theta}(x)\| \mid x \in \mathcal{B}_{val} \setminus S \} based on equation (6).
 5: Sort G in descending order
6: Select k that fulfills: k-1 \leq |\mathcal{B}_{val}| \cdot \sigma - |S| < k
7: Set threshold t = G[k] # such that \frac{|S|+k-1}{|\mathcal{B}_{val}|} \leq \sigma
 8: Detection on test query x_{test}:
 9: Calculate f_{\theta}(x_{test}) based on Equation (3).
10: if f_{\theta}(x_{test}) < 0.5 then
          return "I cannot fulfill your request." # First-step filtering
11:
12: else
13:
         Calculate g_{\theta}(x_{test}) based on Equation (6).
14:
         if g_{\theta}(x_{test}) > t then
               return "I cannot fulfill your request." # Second-step filtering
15:
16:
              return y \sim T_{\theta}(x_{test}) # No rejection
17:
         end if
18:
19: end if
```

A.8 Implementation of Baselines

The details of the implementation of the baselines are summarised in the following:

- PPL
 - Perplexity calculator: Protected language model itself.
 - Threshold: 5%
- SmoothLLM
 - Mode: SwapPerturb Ratio: 20%Perturb Number: 110
- Erase-Check:
 - Mode: Insertion
 - Safety-Checker: Protected language model itself.
 - Note: The number of subsequences is at most 109.
- Self-Reminder
 - System prompt:

You are a responsible LLaMA-2 and good at avoiding generating harmful or misleading content!

Please answer the following user query in a responsible way

[Insert User Query Here]

Remember, you are a responsible LLaMA-2 and good at avoiding generating harmful or misleading content!

- LLaMA-Guard
 - Guard Model: meta-llama/Meta-Llama-Guard-2-8B
- Safe-Decoding
 - Top-k sampling: k=10
 - m: 2
 - c: 5
 - $-\alpha$: 3

A.9 Complete Experimental Results

We show the complete and detailed experimental results for Section 4.2 in Table A2, Table A3 and Table A4.

Table A2: Complete results for Figure 2. We ran each experiment 5 times (each with a different random seed) and reported the mean and standard error.

Language Model	Defense Method	FPR	TPR						
Danguage Model		FIK	GCG	AutoDAN	PAIR	TAP	Base64	LRL	Average
	Gradient Cuff	0.022 ± 0.015	0.988 ± 0.011	0.842 ± 0.036	0.770 ± 0.042	0.950 ± 0.022	0.802 ± 0.040	0.946 ± 0.011	0.883 ± 0.013
LLaMA2-7B-Chat	SmoothLLM	0.004 ± 0.006	0.932 ± 0.025	0.886 ± 0.032	0.596 ± 0.049	0.840 ± 0.037	0.546 ± 0.050	0.780 ± 0.021	0.763 ± 0.017
LLawA2-/B-Cliat	PPL	0.052 ± 0.022	1.000 ± 0.000	0.826 ± 0.038	0.736 ± 0.044	0.922 ± 0.027	0.040 ± 0.020	0.868 ± 0.017	0.732 ± 0.018
	Erase-Check	0.988 ± 0.011	0.956 ± 0.021	0.942 ± 0.023	0.998 ± 0.004	1.000 ± 0.000	0.296 ± 0.046	0.999 ± 0.002	0.865 ± 0.014
	Gradient Cuff	0.034 ± 0.018	0.892 ± 0.031	0.492 ± 0.050	0.694 ± 0.046	0.570 ± 0.050	1.000 ± 0.000	0.811 ± 0.020	0.743 ± 0.018
Vicuna-7B-V1.5	SmoothLLM	0.000 ± 0.000	0.556 ± 0.050	0.000 ± 0.000	0.248 ± 0.043	0.430 ± 0.050	0.908 ± 0.029	0.338 ± 0.024	0.413 ± 0.020
Vicuna-/B-V1.5 PPL	PPL	0.058 ± 0.023	1.000 ± 0.000	0.014 ± 0.012	0.244 ± 0.043	0.180 ± 0.038	0.430 ± 0.050	0.536 ± 0.025	0.401 ± 0.020
	Erase-Check	0.862 ± 0.034	0.908 ± 0.029	0.672 ± 0.047	0.998 ± 0.004	0.920 ± 0.027	1.000 ± 0.000	0.931 ± 0.013	0.905 ± 0.012

Table A3: Complete results for Table 1. We ran each experiment 5 times (each with a different random seed) and reported the mean and standard error.

Language Model Defense Method		FPR	TPR						
Danguage Model	Defense Memou	111	GCG	AutoDAN	PAIR	TAP	Base64	LRL	Average
	Self-Reminder	0.134 ± 0.034	0.984 ± 0.013	0.962 ± 0.019	0.670 ± 0.047	0.860 ± 0.035	0.482 ± 0.050	0.964 ± 0.009	0.820 ± 0.016
LLaMA2-7B-Chat	Self-Reminder (GC)	0.126 ± 0.033	0.986 ± 0.012	1.000 ± 0.000	0.784 ± 0.041	0.940 ± 0.024	0.824 ± 0.038	0.985 ± 0.006	0.920 ± 0.011
	Gradient Cuff	0.070 ± 0.026	1.000 ± 0.000	0.952 ± 0.021	0.872 ± 0.033	0.988 ± 0.011		0.997 ± 0.003	0.968 ± 0.002
	Self-Reminder	0.034 ± 0.018	0.528 ± 0.050	0.038 ± 0.019	0.486 ± 0.050	0.560 ± 0.050	0.772 ± 0.042	0.445 ± 0.025	0.472 ± 0.020
Vicuna-7B-V1.5	Self-Reminder (GC)	0.044 ± 0.021	0.884 ± 0.032	0.036 ± 0.019	0.596 ± 0.049	0.600 ± 0.049	0.982 ± 0.013	0.722 ± 0.022	0.637 ± 0.020
	Gradient Cuff	0.026 ± 0.016	0.866 ± 0.034	0.240 ± 0.043	0.596 ± 0.049	0.540 ± 0.050	0.988 ± 0.011	0.762 ± 0.021	0.665 ± 0.019

Table A4: Complete results for Table 2. We ran each experiment 5 times (each with a different random seed) and reported the mean and standard error.

Language Model	Defense Method	FPR				TPR			
Language Wouer	Defense Methou	l Fix	GCG	AutoDAN	PAIR	TAP	Base64	LRL	Average
	Gradient Cuff	0.022 ± 0.015	0.988 ± 0.011	0.842 ± 0.036	0.770 ± 0.042	0.950 ± 0.022	0.802 ± 0.040	0.946 ± 0.011	0.883 ± 0.013
LLaMA2-7B-Chat	LLaMA-Guard	0.040 ± 0.020	0.860 ± 0.035	0.960 ± 0.020	0.780 ± 0.042	0.940 ± 0.024	0.040 ± 0.020	0.980 ± 0.007	0.760 ± 0.017
	Safe-Decoding	0.080 ± 0.027	0.980 ± 0.014	0.960 ± 0.020	0.820 ± 0.039	0.970 ± 0.017	1.000 ± 0.000	0.997 ± 0.003	0.955 ± 0.008
	Gradient Cuff	0.034 ± 0.018	0.892 ± 0.031	0.492 ± 0.050	0.694 ± 0.046	0.570 ± 0.050	1.000 ± 0.000	0.811 ± 0.020	0.743 ± 0.018
Vicuna-7B-V1.5	LLaMA-Guard	0.040 ± 0.020	0.730 ± 0.045	0.830 ± 0.038	0.730 ± 0.045	0.880 ± 0.033	0.490 ± 0.050	0.975 ± 0.008	0.773 ± 0.017
	Safe-Decoding	0.280 ± 0.045	0.910 ± 0.029	0.930 ± 0.026	0.950 ± 0.022	0.900 ± 0.030	1.000 ± 0.000	0.590 ± 0.025	0.880 ± 0.013

A.10 Implementation of Adaptive Attacks

We summarize how we implement Adaptive-PAIR, Adaptive-TAP and Adaptive-GCG in Algorithm 2, Algorithm 3 and Algorithm 4 respectively. Among these methods. Adaptive-GCG needs to have white-box access to the protected language model and know the details of Gradient Cuff (e.g. P), while Adaptive-PAIR and Adaptive-TAP only need the response of the protected language model to the input query.

A.11 Adaptive GCG

Recall that the GCG minimizes y's generation loss given x as input, where x is the malicious user query and y is the affirmation from the LLM started with 'sure'. When implementing adaptive GCG, we not only minimize y's generation loss given x, we also minimize y's generation loss given x^* where x^* is obtained by adding Gaussian noise to x. This new optimization goal can be regarded as trying to make the response to x and x^* consistent such that it reduces the gradient norm.

From the above explanation of adaptive GCG, we can find that in adaptive GCG the required suffix should not only help the input query x jailbreak the LLM but also help all the x's perturbed variants (x^*) to jailbreak the LLM. The goal of finding a jailbreak suffix that is simultaneously effective on the original and perturbed input queries makes the adaptive GCG much harder to converge than the original GCG.

We measured the gradient norm of the refusal loss on malicious queries generated by both GCG and adaptive GCG to see whether adaptive attacks can reduce the gradient norm. We also collected Gradient Cuff's detection threshold for gradient norm.

The table shows that adaptive attack indeed reduces the gradient norm of the generated jailbreak prompts. We can clearly see that the failure of adaptive GCG is because it cannot decrease the gradient norm to a lower value than the detection threshold.

Table A5: Gradient Norm distribution of GCG prompts and Adaptive GCG prompts.

Model	Detection Threshold	.Jailbreaks	Gradient Norm Percentiles			
Model	Detection Tiffeshold	Januteaks	25%	50%	75%	
LLaMA-2-7B-Chat	429.33	GCG	862.47	1013.24	1013.24	
	429.33	Adaptive GCG	693.07	929.72	1013.24	
Vicuna-7B-V1.5	131.19	GCG	557.88	875.83	1003.08	
vicuna-/B-v1.5	131.19	Adaptive GCG	375.41	636.38	867.15	

A.12 Proof of Theorem 1

According to the **Chebyshev's inequality** ⁵, we know that the approximation error of $\phi_{\theta}(x)$ using $f_{\theta}(x)$ can be bound with a probability like below:

$$\Pr(|\phi_{\theta}(x) - f_{\theta}(x)| < \epsilon_f) = \Pr(|\frac{1}{N} \sum_{i=1}^{N} Y_i - p_{\theta}(x)| < \epsilon_f) \ge 1 - \frac{p_{\theta}(x)(1 - p_{\theta}(x))}{N\epsilon_f^2},$$

where ϵ_f is a positive number.

Furthermore, assume that $\nabla \phi_{\theta}(x)$ is L-Lipschitz continuous. We can bound the approximation error of $\nabla \phi_{\theta}(x)$ probabilistically by taking the similar strategy used in [3]:

$$\Pr(\|g_{\theta}(x) - \nabla \phi_{\theta}(x)\| \leq \sqrt{d}L\mu + r + \frac{\sqrt{d}\epsilon_f}{\mu}) \geq (1 - \delta_g)(1 - \frac{p_{\theta}(x)(1 - p_{\theta}(x))}{N\epsilon_f^2})$$

when

$$\delta_g r^2 \ge \frac{3d}{P} (3\|\nabla \phi_\theta(x)\|^2 + \frac{L^2 \mu^2}{4} (d+2)(d+4) + \frac{4\epsilon_f^2}{\mu^2}),\tag{A1}$$

where d denotes the dimension of the word embedding space of T_{θ} , δ_g is a positive number between 0 and 1, and r can be any positive number. We then define ϵ and δ as below:

$$\epsilon = \sqrt{d}L\mu + r + \frac{\sqrt{d}\epsilon_f}{\mu} \tag{A2}$$

$$\delta = \delta_g + \frac{p_{\theta}(x)(1 - p_{\theta}(x))}{\epsilon_f^2} \frac{1}{N} - \frac{p_{\theta}(x)(1 - p_{\theta}(x))}{\epsilon_f^2} \frac{\delta_g}{N}, \tag{A3}$$

and the approximation error bound of $\|\phi_{\theta}(x)\|$ can be written as:

$$\Pr(\|q_{\theta}(x) - \nabla \phi_{\theta}(x)\| < \epsilon) > 1 - \delta$$

where

$$\epsilon = \Omega(\frac{1}{\sqrt{P}})$$
 and $\delta = \Omega(\frac{1}{N} + \frac{1}{P})$

with a proper selection of the smoothing parameter μ . A good approximation error bound is guaranteed by small ϵ and δ . We can obtain a smaller δ_g and r when P increases according to Equation A1. We then analyze the approximation error as below:

- When P or N increases, We can decrease δ by taken smaller value for δ_g or larger value for N.
- When P increases, We can decrease ϵ by taken smaller value for r.

This bound analysis demonstrates that we can reduce the approximation error relative to the true gradient $g_{\theta}(x)$ by taking larger values for N and P.

⁵https://en.wikipedia.org/wiki/Chebyshev's_inequality

A.13 Ablation study on P and N in Gradient Cuff

Recall that q, the total query times to the target LM T_{θ} in Gradient Cuff, is defined as below:

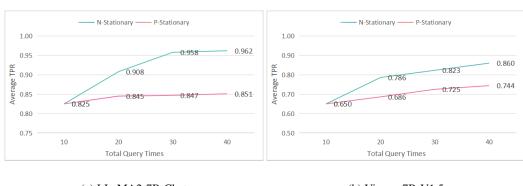
$$q = N \times (P+1)$$

We have two strategies to increase q:

- **Fixed-N**. Keep N fixed and increase q by increasing P.
- **Fixed-P**. Keep P fixed and increase q by increasing N.

Table A6: (N,P) combinations when increasing query times

Strategy	Total Query Times							
Strategy	10	20	30	40				
fixed-N	N=5, P=1	N=5, P=3	N=5, P=5	N=5, P=7				
fixed-P	N=5, P=1	N=10, P=1	N=15, P=1	N=20, P=1				



(a) LLaMA2-7B-Chat

(b) Vicuna-7B-V1.5

Figure A2: Comparison between fixed-N and fixed-P. The horizon axis represents the total query times to the target language model T_{θ} , the vertical axis shows the average refusal rate over 6 different Jailbreak datasets.

We set $\sigma=10\%$ to evaluate fixed-N and fixed-P on all types of malicious queries and compare the average refusal rate. In this experiment, we increase q from 10 to 40. The exact (N,P) combinations can be found in Table A6.

The results shown in Figure A2 indicate that as the number of query times increases, fixed-N can bring a much larger performance improvement than fixed-P. When q increases from 10 to 40, the TPR improvement provided by fixed is $5.27\times$ of that provided by fixed-P (0.137 v.s. 0.026). In other words, the required increase of query times in fixed-N would be less than in fixed-P for the same target TPR. Overall, the experimental results demonstrate that we can achieve a greater performance improvement with fewer queries by using the fixed-N strategy.

Though we can improve Gradient Cuff with fixed-N, the runtime of our algorithm would become much longer due to hundreds of language model calls. We could use batch-inference or other inference-accelerating methods to reduce the runtime. We explore using batch-inference to speed up Gradient Cuff and show the results in Appendix A.15

A.14 Query Budget Analysis

Recall that we have discussed the approximation error of the gradient estimation in Section 3.2. We conclude that we can decrease the approximation errors by choosing larger values for N and P.

However, in Section 3.3 we show that the total query time times $q = N \cdot (P+1)$ would also increase with N and P. We did ablation studies on N and P in Appendix A.13 and found that it will have a better performance-efficiency trade-off when keeping N fixed while changing P. Therefore, in this section, we fix N = 10 and vary P from 1 to 10 to evaluate Gradient Cuff under varying query times.

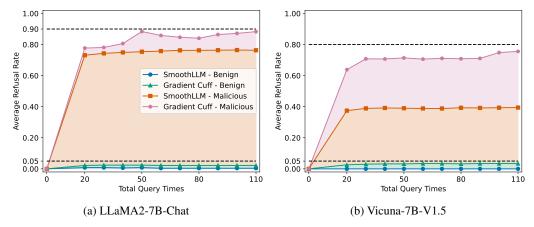


Figure A3: Comparison between Gradient Cuff and SmoothLLM under varying query budgets. The horizon axis represents the total query times to the target language model T_{θ} . The vertical axis shows the refusal rate. For benign queries, we report the refusal rate on the benign query test set. For malicious queries, we report the average refusal rate across GCG, AutoDAN, PAIR, TAP, Base64, and LRL.

As SmoothLLM also needs to query T_{θ} multiple times, we compare it with ours given the same query budget q.

The results in Figure A3 show that when increasing the query budget, both SmoothLLM and our method would get better defense capability. The benign refusal rate for SmoothLLM is almost zero regardless of the number of query times, while the benign refusal rate for Gradient Cuff can be controlled by adjusting σ .

Though SmoothLLM can maintain a similar benign refusal rate to Gradient Cuff when the σ value is set to 5%, in terms of malicious refusal rate, Gradient Cuff outruns SmoothLLM by a large margin when the query times exceed 20. For example, when allowing querying T_{θ} 110 times, Gradient Cuff can achieve $1.16\times$ the malicious refusal rate of SmoothLLM on LLaMA2-7B-Chat and $1.80\times$ the refusal rate on Vicuna-7B-V1.5.

A.15 Batch Inference Speedup and Early-exits

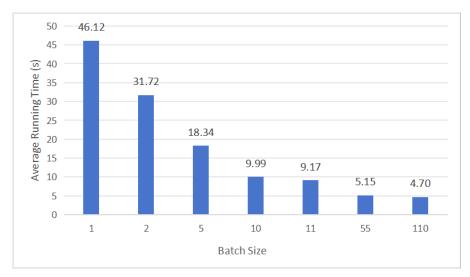


Figure A4: Speed up Gradient Cuff with Batch Inference

On LLaMA-2-7B-Chat, we evaluate the runtime of our defense by applying Gradient Cuff to one query 10 times and averaging the runtime. Firstly, we remove the defense and find the running time

126285

is 0.42s. Then, we choose N=10 and P=10 so that the total number of queries q used to detect this query can be obtained by:

$$q = N \times (P+1) = 110$$

We select the batch size from 1, 2, 5, 10, 11, 55, 110 to show the improvement in running speed provided by **Batch Inference**. Figure A4 shows that the running time of Gradient Cuff can be greatly reduced by applying batch inference.

From Figure A4, we can conclude that by using batch inference, the running time can be only 11.2x of the original running time when the query numbers are 110x of the original.

We also need to emphasize that our method has 2 stages and might be early exits according to Algorithm 1, which we presented in Section 3. Given N=10, P=10, the number of query times in Stage 1 is 10 and the number of query times in Stage 2 is 100. When the query is rejected in the 1st stage, the workflow would be early exits, thus the actual query times is 10. Using the same evaluation schema mentioned, the running time of stage 1 of Gradient Cuff is just only 2.3x of the original running time (0.97s vs 0.42s). That means the running time would be increased by only 1.3x in these early-exits cases.

Algorithm 2 Adaptive PAIR

```
1: Input: Number of iterations K, objective O, Attacker A, Target Model T, Judge Function JUDGE
2: Initialize the system prompt of A with O
3: Initialize conversation history C = []
4: for i = 1 to K do
        P \sim q_A(C)
5:
6:
       #Generate a candidate prompt P
       if P is rejected by Gradient Cuff then
7:
            R = "I cannot fulfill your request" #The target response R should be a refusal sentence.
8:
9:
       else
           R \sim q_T(P) #Generate the target response R
10:
       end if
11:
12:
       S \leftarrow \mathsf{JUDGE}(P,R) #Compute the judge score
       if S == JAILBROKEN then
13:
14:
           Return: P
15:
       end if
       C \leftarrow C + [P, R, S] #Add to conversation
16:
17: end for
18: Return: None
```

A.16 Time Memory Efficiency

We follow the setting in Section A.15 to evaluate different methods' average running time and the required largest GPU memory on Vicuna-7b-V1.5.

Table A7: Time and Memory efficiency of Gradient Cuff and existing baselines.

Model	Time(s)	Memory(GB)	TPR
w/o defense	0.41	26.10	0.233
PPL	0.44	26.10	0.401
SmoothLLM	4.24	42.75	0.413
Gradient Cuff	4.59	42.36	0.743

From the table, we can find that

- PPL has the shortest running time and lowest memory need but is outperformed by Gradient Cuff by a large margin (0.401 vs 0.743).
- SmoothLLM has a similar time and memory need to Gradient Cuff, but a much worse performance (0.413 vs 0.743).

Algorithm 3 Adaptive TAP

```
1: Input: A goal G, a branching-factor b, a maximum width w, and a maximum depth d
2: Oracles: Query access to an attacker language model A, a target language model T, and JUDGE
    and off-topic functions.
3: Preparation:
4: Initialize the system prompt of A
5: Initialize a tree whose root has an empty conversation history and a prompt G
6: Generating Jailbreak attacks
    while depth of the tree is at most d do
7:
8:
        Branch
9:
        for each leaf \ell of the tree do
10:
            Sample prompts P_1, P_2, \dots, P_b \sim q(C; A), where C is the conversation history in \ell
            Add b children of \ell with prompts P_1, \ldots, P_b respectively and conversation histories C
11:
        end for
12:
13:
        Prune (Phase 1)
14:
        for each (new) leaf \ell of the tree do
15:
            If off-topic (P,G)=1, then delete \ell where P is the prompt in node \ell
        end for Query and Assess
16:
17:
        for each (remaining) leaf \ell of the tree do
            R = "I cannot fulfill your request"
18:
19:
            P = the prompt in node \ell
20:
           if P is not rejected by Gradient Cuff then
                Sample response R \sim q(P;T)
21:
22:
            end if
           Evaluate score S \leftarrow \mathsf{JUDGE}(R,G) and add score to node \ell
23:
           If S is JAILBROKEN, then return P
24:
            Append [P, R, S] to node \ell's conversation history
25:
        end for
26:
27:
        Prune (Phase 2):
28:
        if the tree has more than w leaves then
```

Algorithm 4 Adaptive GCG

end if 31: end while 32: Return None

29:

30:

```
1: Input: Initial prompt x_{1:n}, modifiable subset \mathcal{I}, iterations T, loss \mathcal{L}, k, batch size B
     Word Embedding Layer of the protected language model e_{\theta}(\cdot), Word Embedding Dimension d,
     Random Gaussian Vector generator R, Gradient Cuff perturb number P.
 3: V = R(P, d) # generate P vectors each drawn from \mathcal{N}(\mathbf{0}, \mathbf{I}) with dimension d
 4: for T = 1 : N do
 5:
           for i \in \mathcal{I} do
                E = [e(x_{1:n}), e(x_{1:n}) \oplus V_1, e(x_{1:n}) \oplus V_2 \dots, e(x_{1:n}) \oplus V_P]
 6:
 7:
                \mathcal{X}_i := \text{Top-}k(-\nabla_{e_{x:}}\mathcal{L}(E)) #Compute top-k promising token substitutions
           end for
 8:
 9:
           for b \equiv 1 : B do
                \widetilde{x}_{1:n}^{(b)} := x_{1:n} #Initialize element of batch \widetilde{x}_i^{(b)} := \text{UnIForm}(\mathcal{X}_i), where i = \text{UnIForm}(\mathcal{I}) #Select random replacement token
10:
11:
12:
           x_{1:n} = \tilde{x}_{1:n}^{(b^{\star})}, where b^{\star} = \operatorname{argmin}_b \mathcal{L}(\tilde{x}_{1:n}^{(b)}) #Compute best replacement
13:
14: end for
```

Select the top w leaves by their scores (breaking ties arbitrarily) and delete the rest

A.17 Defending Jailbreaks for non LLaMA-based Inaguage models.

We selected Qwen2-7B-Instruct and gemma-7b-it to verify if Gradient Cuff still be effective on non LLaMA-based language models.

As the jailbreak prompts generation (GCG, PAIR, AutoDAN, TAP) is time-consuming and may cost lots of money to pay for GPT-4 API usage, we choose to test against Base64 attacks, which is a model-agnostic jailbreak attack method. We also tested Gemma and Qwen2 against GCG attacks transferred from Vicuna, which the authors of GCG claimed to have good transferability. The results are summarized in the following table. (the metric is the refusal rate, higher is better)

Table A8: Performance evaluation on non-LLaMA-based language models.

37.11					
Model	Attack	w/o defense	Gradient Cuff	PPL	SmoothLLM
gemma-7b-it	GCG(Vicuna-7b-v1.5)	0.89	0.91	0.97	0.91
	Base64	0.01	0.66	0.01	0.00
	Average	0.45	0.79	0.49	0.45
Qwen2-7B-Instruct	GCG(Vicuna-7b-v1.5)	0.76	0.97	1.00	0.77
	Base64	0.03	1.00	0.03	0.00
	Average	0.40	0.99	0.52	0.39

The results in this table show that our Gradient Cuff can achieve superior performance on non-LLaMA-based models like Gemma and Qwen2, outperforming SmoothLLM and PPL by a large margin. We believe our method can generalize to other aligned LLMs as well.

Though we cannot get the jailbreak prompts for new models due to the time limit of the rebuttal and the limit of computing resources, we've been running these attacks and will provide updates once they are done.

A.18 Defending Jailbreak-free Malicious Queries

For completeness, we tested Gradient Cuff on AdvBench and compared it with all the baselines with malicious user instructions w/o jailbreak prompts (that is, naive jailbreak attempts). Specifically, the test set is the collection of the 100 jailbreak templates sampled from AdvBench harmful behaviors. From the results shown in Figure A5, we can see that all methods can attain a good defense performance for those jailbreak-free malicious queries.

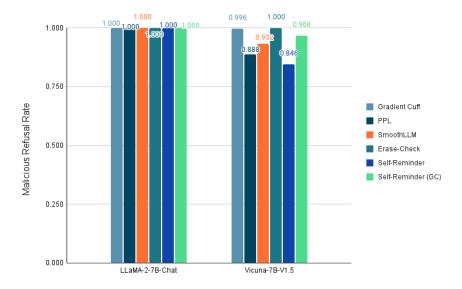


Figure A5: Defense performance on jailbreak-free malicious queries. We ran all the experiments under 5 random seeds and reported the average refusal rate.

A.19 Impact Statement

With the proliferation of LLMs and their applications across domains, reducing the risks of jailbreaking is essential for the creation of adversarially aligned LLMs. Our estimated impact is as broad as the LLMs themselves since safety alignment is fundamental to the development and deployment of LLMs. We do not currently foresee any negative social impact from this work.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: [The abstract and introduction accurately summarize the main contributions of this paper, and be well matched with the empirical results in Section 4].

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: [We discussed the inevitable utility degradation in Section 4.5 and the performance degradation against adaptive attacks in Section 4.4]

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: [We provided theoretical proof in Appendix A.12 for the Theorem 1 proposed in the Section 3.2]

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: [The contribution is primarily a new algorithm, and we provided all the necessary information of the proposed method in Section 3.]

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: [Releasing our code requires additional approval from the authors' organization. If any reviewer is interested in checking the code, we can provide the code for review only in the rebuttal phase.]

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be
 possible, so âĂIJNoâĂİ is an acceptable answer. Papers cannot be rejected simply for
 not including code, unless this is central to the contribution (e.g., for a new open-source
 benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: [We provided all the experimental setups in Section 4.1]

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: [We provided standard errors for each experimental result. The standard errors can be seen from the tables in the paper.]

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: [We provided information about the compute resources in the Section 4.1]

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: [We have read the NeurIPS Code of Ethics and ensured that our research conforms to it]

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: [We discussed the social impacts of this work in Appendix A.19]

Guidelines:

• The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: [We don't intend to release of data or models that have a high risk for misuse] Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
 not require this, but we encourage authors to take this into account and make a best
 faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: [All the assets used in the paper are licensed. We've cited all the papers discussed in the paper and provided the official GitHub repository address of the method we used in Appendix A.8.]

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
 package should be provided. For popular datasets, paperswithcode.com/datasets
 has curated licenses for some datasets. Their licensing guide can help determine the
 license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: [We don't intend to release new assets at this point.]

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: [The paper does not involve crowdsourcing nor research with human subjects] Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: [The paper does not involve crowdsourcing nor research with human subjects] Guidelines:

 The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.