# Beyond Surface Alignment: Rebuilding LLMs Safety Mechanism via Probabilistically Ablating Refusal Direction

**Anonymous EMNLP submission**

## Abstract

Jailbreak attacks pose persistent threats to large language models (LLMs). Current safety alignment methods have attempted to address these issues, but they experience two significant limitations: insufficient safety alignment depth and unrobust internal defense mechanisms. These limitations make them vulnerable to adversarial attacks such as prefilling and refusal direction manipulation. We introduce DeepRefusal, a robust safety alignment framework that overcomes these issues. DeepRefusal forces the model to dynamically rebuild its refusal mechanisms from jailbreak states. This is achieved by probabilistically ablating the refusal direction across layers and token depths during fine-tuning. Our method not only defends against prefilling and refusal direction attacks but also demonstrates strong resilience against other unseen jailbreak strategies. Extensive evaluations on four open-source LLM families and six representative attacks show that DeepRefusal reduces attack success rates by approximately 95%, while maintaining model capabilities with minimal performance degradation. Our code is available at: anonymous git.

## 1 Introduction

Large Language Models (LLMs) (Achiam et al., 2023; Mistral AI, 2024; Team et al., 2024; Grattafiori et al., 2024) have demonstrated impressive performance across various natural language tasks, but ensuring them behave safely and reliably remains a significant challenge. Saftey alignment efforts, such as refusal training, have endowed LLMs with the ability to provide refusal responses to inappropriate and toxic prompts, as shown in Figure 1a. However, these safeguards frequently fail when confronted with adversarial Jailbreak attacks that disguise harmful requests to induce LLMs to circumvent refusal behaviors and trap in Jailbreak states (Wei et al., 2023a; Zou et al., 2023b; Chao et al., 2023), as shown in Figure 1b.

At present, numerous defense strategies have been proposed to defend against jailbreak attacks. Adversarial training (Mazeika et al., 2024) tries to enhance refusal robustness of LLMs against such attacks by training on jailbreak samples. Circuit-Breaker (Zou et al., 2024) suppresses harmful activations within the model's hidden states (Figure 1c) to bolster resilience. However, our experiments show that existing defense strategies, even state-of-the-art ones, can be easily bypassed by prefilling attacks and refusal direction attacks, and are susceptible to transfer attacks. Specifically, we find that existing defenses remain predominantly surface-level and exhibit two significant limitations. (1) **Insufficient Safety Alignment Depth**: Current safety alignment methods focus on suppressing toxicity within the initial few response tokens, while overlooking the harmfulness of subsequent tokens. As a result, the inherent safeguard of LLMs can be bypassed by manipulating the initial response tokens (i.e, prefilling attack) (Qi et al., 2025; Zhou et al., 2024a; Vega et al., 2023). (2) **Unrobust Internal Defense Mechanisms**: The shallow internal defense mechanisms demonstrate limited resilience to advancing jailbreak methodologies. Attackers can iteratively refine input prompts to jailbreak LLMs. Additionally, simple techniques such as refusal direction ablation (Arditi et al., 2024) can readily circumvent defenses[1].

To mitigate the above limitations inherent in current safety alignment methods, we present **DeepRefusal**, an innovative framework designed for deep and robust safety alignment. Unlike traditional approaches that focus merely on surface-level fine-tuning, DeepRefusal introduces adversarial pressure directly into the model's representation space during the fine-tuning process, as illustrated in Figure 1d. By probabilistically ablating refusal di-

---

[1] Jailbreaking by fine-tuning is **not** within the scope of this paper. Here are some orthogonal work: (Qi et al., 2023; Tamirisa et al., 2024; Huang et al., 2024).
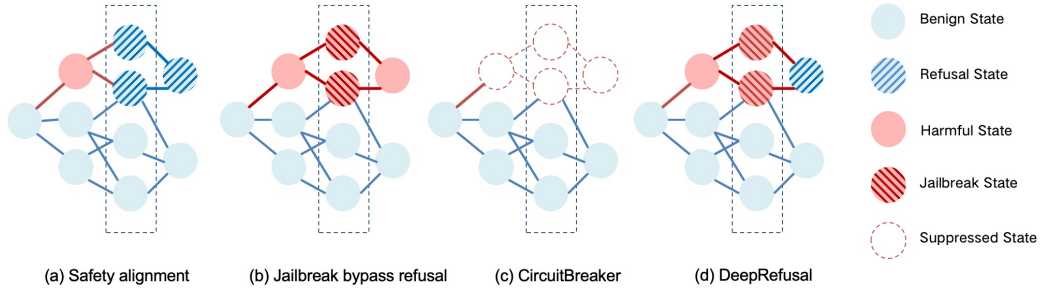
Figure 1: Conceptual diagram of jailbreak and safety mechanisms at the representation level: (a) Safety alignment prevents the model from answering harmful requests by activating the concept of refusal. (b) Jailbreak attacks bypass the refusal behavior of the model through carefully designed prompts. (c) CircuitBreaker (Zou et al., 2024) suppresses harmfulness within the hidden states (still vulnerable to attacks targeting token depth or internal refusal mechanisms). (d) DeepRefusal simulates jailbreak scenario at the representation level, forcing the model to reactivate the refusal behavior, which fundamentally enhances the robustness of the model's refusal mechanism.

rections across multiple layers and token depths, DeepRefusal effectively simulates potential jailbreak scenarios internally. The refusal direction is a direction in the model's internal representation space that is strongly associated with producing refusal responses. This unique approach compels the model to develop a more robust refusal mechanism, ensuring it can resist harmful outputs at every token depth and layer depth, thereby enhancing its safety alignment capabilities significantly.

Our contributions are summarized as follows:

- We identify two key technical challenges in robust safety alignment and show that attacks along the refusal direction are a transferable, under-addressed threat to state-of-the-art safety alignment methods.

- We design DeepRefusal, a framework that trains LLMs to rebuild safety mechanisms from jailbreak states by simulating adversarial conditions across various token and layer depths, bridging the gap between surface-level alignment and robust internal defense

- Through extensive experiments, we demonstrate that DeepRefusal significantly improve alignment robustness, reducing attack success rates by approximately 95%, while effectively defending against prefilling attacks, refusal direction attacks, and showing strong robustness to other unseen jailbreak attacks.

## 2 Related Works

**Jailbreak Attacks on LLMs**. While LLMs demonstrate impressive performance, they remain vulnerable to adversarial inputs that bypass safety constraints. Early works explored manually crafted prompts to elicit harmful outputs (Wei et al., 2023a; Mazeika et al., 2024), while more recent approaches have employed automated methods. GCG attacks generate adversarial suffixes via gradient-based optimization (Zou et al., 2023b), and CodeAttack exploits insufficient alignment in code domains (Ren et al., 2024). AutoDAN employs genetic operations to evolve jailbreak prompts (Liu et al., 2023), and PAIR uses attacker–victim loops to iteratively refine jailbreak prompts (Chao et al., 2023). Prefilling attacks manipulate early token activations to influence model outputs (Vega et al., 2023). Recent studies have revealed that refusal direction (Arditi et al., 2024) can be used for jailbreak. Our experiments confirm that these refusal direction attacks, achieve high success rates and exhibit strong transferability within fine-tuned models. Furthermore, we find that current state-of-the-art defense techniques fail to defend against such attacks, highlighting the need for representation-level robustness.

**LLM Safety Alignment**. Research into LLM safety alignment has yielded several defensive strategies against adversarial attacks. For instance, R2D2 (Mazeika et al., 2024) and CAT (Xhonneux et al., 2024) utilize adversarial training, incorporating optimized prompts or input-level perturbations to strengthen model resilience. Meanwhile, Latent Adversarial Training (LAT) (Sheshadri et al., 2024; Casper et al., 2024) focuses on enhancing the model's resistance to harmful outputs by targeting its internal residual stream. However, these approaches often demand substantial computational

resources. Additionally, a method proposed in (Qi et al., 2025) aims to achieve deep safety alignment via data augmentation, yet it proves less effective against previously unseen attacks. CircuitBreaker (Zou et al., 2024), which defends against jailbreak attacks by suppressing harmful hidden states, tends to produce gibberish generation. This drawback significantly limits its practical utility in applications where clear and timely refusal responses are essential. responses are preferred.

Besides the above works, some solutions focus on designing input/output filters to screen potentially harmful content (Alon and Kamfonas, 2023; Inan et al., 2023) or intervening during model decoding at inference-time (Robey et al., 2024; Xu et al., 2024a). However, these methods do not fundamentally enhance model alignment. Instead, they necessitate additional components, thereby introducing extra overhead during testing. It is important to note that these approaches are orthogonal to safety fine-tuning, and can be used in conjunction during test-time. Thus, we do not compare our methods with these in this paper.

**Representation Engineering and Refusal Direction**. Recently, analysis and manipulation techniques for the internal representations (i.e., activations) of LLMs have gained widespread attention. (Zou et al., 2023a) formally introduced Representation Engineering (RepE), which draws on insights from cognitive neuroscience to improve the transparency of AI systems. Furthermore, (Arditi et al., 2024) demonstrated that refusal in LLMs is mediated by a single direction, termed the refusal direction, and elucidated the mechanism of jailbreaking at the representation level. Additionally, (Panickssery et al., 2023; Stickland et al., 2024) showcased how representation-level interventions can effectively control model output behavior. Our approach leverages these insights into the explainability of jailbreaking and LLM safety mechanisms to enhance model alignment and robustness.

## 3 Observations and Motivation

### 3.1 Current Safety Fine-Tuning Made on the First Few Tokens

The majority of current safety alignment methods suffer from a shared limitation: they over-rely on aligning the model's generation within the initial few tokens of the response (Yuan et al., 2024; Qi et al., 2025). This gives rise to an inherent vulnerability, termed the "alignment depth problem". This
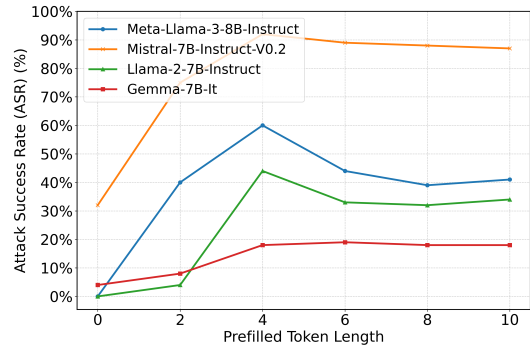


Figure 2: Attack success rate with different lengths of harmful tokens prefilled.

problem is rooted in a widely held yet unrealistic assumption: when encountering harmful prompts, a well-aligned model should instantaneously refuse, i.e., produce refusal tokens like "I'm sorry..." at the very beginning of its output. However, recent work (Wei et al., 2023b; Zhou et al., 2024b) have shown that once this initial safeguard is bypassed, the model often freely generates harmful content.

This behavior is exploited by nearly all recently proposed jailbreak attacks. For instance, DeRTa (Yuan et al., 2024) finds that refusal responses typically place the first refusal token within the first 20 output tokens. If a prompt succeeds in suppressing refusal behavior in this early region, the likelihood of a successful jailbreak increases dramatically. Similarly, Qi et al. (2025) show through KL divergence analysis that safety fine-tuning disproportionately aligns early tokens, leaving later tokens insufficiently aligned.

From the attacker's perspective, Wei et al. (2023b) argue that jailbreaks are essentially about suppressing early refusal. This is further exemplified by the prefilling attack (Vega et al., 2023), which prepends a harmful context before the actual prompt to mislead the model's early generation. To validate this, we conduct a prefilling attack experiment shown in Figure 2. We randomly select 100 harmful instructions from AdvBench (Zou et al., 2023b), and prepend prefilled harmful tokens of increasing length. The results demonstrate a clear trend: as the length of the prefilled content increases, the attack success rate rises significantly.

The current evidence collectively indicates that merely relying on shallow alignment at the beginning of a response is far from adequate for ensuring robust safety. This underscores the critical vulnerability arising from the alignment depth problem and highlights the necessity for deeper safety align-

3

ment mechanisms that extend beyond the initial token positions. Insteadly, they should be deeply embedded within the model's internal representations, consistently spanning all stages of response generation.

## 3.2 Jailbreak Inhibits the Activation of Refusal Direction

Recent works (Arditi et al., 2024; Zou et al., 2023a; Xu et al., 2024b) have demonstrated that LLMs internally encode a refusal direction, which is a distinct activation pattern corresponding to refusal generation. To extract this direction, we adopt the methodology from Arditi et al. (2024). Formally, for each layer $l \in [L]$ and post-instruction token position $i \in [I]$, we calculate the mean activation $\mu_i^{(l)}$ for harmful prompts from $D_{\text{harmful}}^{(\text{train})}$ and $\nu_i^{(l)}$ for harmless prompts from $D_{\text{benign}}^{(\text{train})}$:

$$\mu_i^{(l)} = \frac{1}{|D_{\text{harmful}}^{(\text{train})}|} \sum_{t \in D_{\text{harmful}}^{(\text{train})}} h_i^{(l)}(t) \quad (1)$$

$$\nu_i^{(l)} = \frac{1}{|D_{\text{benign}}^{(\text{train})}|} \sum_{t \in D_{\text{benign}}^{(\text{train})}} h_i^{(l)}(t) \quad (2)$$

where $h_i^{(l)}(t)$ denotes the hidden activation vector at layer $l$ and position $i$ for input $t$. The difference between these two mean vectors defines the candidate refusal direction:

$$r_i^{(l)} = \mu_i^{(l)} - \nu_i^{(l)} \quad (3)$$

This process yields an array of $|I| \times L$ candidate vectors. In line with Arditi et al. (2024), we refine this array through heuristic filtering to identify the most effective single direction $r_{i*}^{(l^*)}$, judging by its capacity to either induce or suppress refusal behavior upon manipulation. Importantly, this analysis reveals that jailbreak prompts actively **inhibit** the model's internal refusal direction. As depicted in Figure 3, when processing adversarial prompts, the cosine similarity between model activations and the refusal direction shows a marked decline. This indicates that jailbreaks are not merely prompt tricks; instead, they **directly alter the model's internal safety-related activations**, effectively bypassing its built-in refusal mechanism.

This observation yields a crucial insight: We can simulate the worst-case impact of a jailbreak attack by **directly ablating** the refusal direction from the LLM's activations. Unlike traditional training
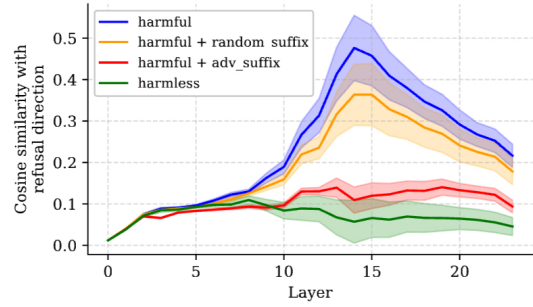


Figure 3: Cosine similarity between the residual stream activation of the last token and refusal direction(Arditi et al., 2024).

methods that rely on specific jailbreak prompts, this activation-level approximation zeroes in on the core mechanism exploited by jailbreak attacks. It removes the necessity for detailed knowledge of specific attack methodologies, making the approach more universal and controllable. More importantly, it enables models to rebuild refusal behaviors even under the extreme internal conditions. This innovation also eliminates the need for a large-scale jailbreak corpus, making the training process more efficient and practical.

## 4 Our DeepRefusal Method

Drawing on these observations, we introduce **DeepRefusal**, a novel approach designed to achieve deep and robust safety alignment for LLMs. DeepRefusal simulates jailbreak scenarios within the model's internal structure without undermining its language modeling capabilities or utility. By doing so, it compels the model to rebuild and reinforce its refusal mechanism across layer depth and token depth. Our goal is to expose the model to these simulated jailbreak scenarios during training. This ensures the model learns to maintain robust refusal behaviors even under adversarial conditions.

In practice, DeepRefusal utilizes two main strategies to accomplish this goal: Layer-wise and Tokewise Probabilistic Activation Ablation (PAA). The former involves probabilistically removing activation patterns across different layers of the model, and the latter focuses on selectively ablating activations at specific token positions within the model's output sequence. By integrating these strategies, we ensure the model encounters a diverse array of simulated adversarial attacks. This all-encompassing method compels the model to develop robust defenses against various jailbreak techniques, thereby maintaining safe and reliable

4

performance even in challenging scenarios.

## 4.1 Obtain the Refusal Direction

A critical aspect of DeepRefusal is obtaining a direction that truly represents the refusal concept. The direction must induce refusal behavior across both benign and harmful samples: when added, it encourages the model to refuse, and when ablated, it prevents the model from refusing. Given different token positions and layers within an LLM may offer a large number of potential directions, we use a heuristic filtering approach to identify the most representative refusal direction (Arditi et al., 2024).

This involves evaluating candidate directions based on their ability to consistently trigger or suppress refusal responses across a diverse set of prompts. Specifically, we define the refusal direction as the vector $r$ that satisfies two key criteria:

**(1) Addition Constraint:** Adding this direction should consistently trigger a refusal response, even for harmless prompts.

$$h^{(l)\prime} \leftarrow h^{(l)} + r^{(l)} \qquad (4)$$

**(2) Ablation Constraint:** Removing this direction should bypass the refusal mechanism, allowing the model to respond to harmful prompts.

$$h' \leftarrow h - \hat{r}\hat{r}^\top h \qquad (5)$$

From a set of candidate directions $\{r_i^{(l)}\}$, we select the single direction $\{r_{i*}^{(l*)}\}$ (denoted $\hat{r}$) that best satisfies the both constraints across validation data. In addition, we tried to dynamically change the refusal direction during fine-tuning, but we found that the direction obtained in this way was very unstable and it was difficult to meet the above two constraints in real time. We finally chose to pre-calculate the optimal single refusal direction offline. And this choice ensures that the direction of ablation is definitely the refusal direction, reducing side-effects on other concept representations and language modeling. Our experiments proved that using the offline refusal direction is sufficient to improve alignment.

## 4.2 Simulate Jailbreak with Layer-wise PAA

To simulate jailbreak scenarios within the model's internal hidden states and prompt the model to re-build its refusal mechanism (Figure 1d), we implement an activation intervention at each of the $L$ layers independently with probability $p$. Specifically, we denote

$$Q_l \sim \text{Bernoulli}(p), \quad l = 1, 2, \ldots, L, \qquad (6)$$

as the indicator that layer $l$ is intervened. Here, $Q_i = 1$ indicates that we apply directional activation ablation at layer $l$. Consequently, Equation 5 can be reformulated as:

$$h' \leftarrow h - Q(\hat{r}\hat{r}^\top h) \qquad (7)$$

Layer-wise PAA allows us to probabilistically remove activation patterns across different layers of the model. This simulates adversarial conditions where the model's internal safety mechanisms are compromised at various layers. By doing so, the model is forced to reinforce its refusal behavior across its entire layer depth. This process significantly enhances the model's robustness against potential jailbreak attacks.

## 4.3 Simulate Jailbreak with Token-wise PAA

Let $x = (x_1, \ldots, x_{|x|})$ denotes the user prompt and $y = (y_1, \ldots, y_{|y|})$ denotes the models's output. The full sequence length is $T = |x| + |y|$. We denote $\pi_\theta(y \mid x)$ as the model's conditional distribution over the response given the prompt. By the chain rule, this can be expanded as:

$$\pi_\theta(y \mid x) = \prod_{t=1}^{|y|} \pi_\theta(y_t \mid x, y_{<t}) \qquad (8)$$

We introduce $M_{l,t} \sim \text{Bernoulli}(p)$ as an indicator variable that determines whether to ablate the refusal direction $\hat{r}$ at layer $l$ and position $t$. During the forward pass, the intervened hidden state is updates as follows:

$$h_t^{(l)\prime} = h_t^{(l)} - Q_l \times M_{l,t} \left( \hat{r}\,\hat{r}^\top h_t^{(l)} \right) \qquad (9)$$

The loss function is then defined as:

$$\mathcal{L} = -\sum_{t=1}^{|y|} \log \pi_\theta\left(y_t \mid x, y_{<t}; \{h_t^{(l)\prime}\}_{l=1}^L\right) \qquad (10)$$

Token-wise PAA selectively targets specific token positions within the model's output sequence. This strategy directly addresses the model's vulnerability to generate harmful content when certain tokens are either manipulated or suppressed. By focusing on these critical positions, the model is trained to recognize and counteract such adversarial manipulations. As a result, it learns to consistently maintain robust refusal behaviors throughout the entire token generation process, thereby significantly enhancing its safety and reliability.

## 4.4 Training Procedure of DeepRefusal

Our final design draws from Qi et al. (2025) and utilizes quadruples in the form of $(x, x', y, y')$, where $x$ denotes benign instructions paired with its corresponding safe response $y$, and $x'$ denotes harmful instructions associated with harmful response $y'$. A pair $(x', y)$ signifies a harmful instruction met with a refusal response, which is the desired result of safe alignment. Conversely, a pair $(x', y')$, typically indicates that the LLM has been jailbroken. Our primary objective is that for any $x'$, the model consistently generates safe response $y$.

This consistency must be maintained even under challenging conditions such as prefilling attacks, where the model encounters modified input sequences $(x', y'_{<t})$. In such scenarios, we aim for the model to persist in generating the safe response $y$ rather than succumbing to the attack and producing the harmful response $y'$. Similarly, in the face of representation-level attacks, such as those described in Eq.(9), which involve ablating the refusal direction, our goal is for the model to reconstruct and maintain the refusal response $y$.

We employ Equation 11 to simulate the prefilling attack, where a harmful prefix is prepended to the response.

$$\pi_\theta(y|x, y'_{\le k}), \quad k \sim \mathrm{Uniform}[20, 25] \qquad (11)$$

To fine-tune the model against such attacks, we use the following objective function:

$$\begin{aligned}\min_\theta \ &\alpha \times \mathbb{E}\big[-\log \pi_\theta(y|x', y'_{\le k}; \{h^{(l)\prime}\}_{l=1}^L)\big] \\ &+(1-\alpha) \times \mathbb{E}\big[-\log \pi_\theta(y|x; \{h^{(l)\prime}\}_{l=1}^L))\big]\end{aligned} \qquad (12)$$

Consistent with Qi et al. (2025), we set $\alpha$ to 0.2. Our training dataset comprises three types of pairs: $(x, y)$ to preserve the model's utility, $(x', y)$ for refusal training, and the augmented pair $(x', y'_{<t}, y)$ to enhance robustness against prefilling attacks.

Algorithm 1 outlines the entire procedure of DeepRefusal. It first obtains a global refusal direction $\hat{r}$, following Arditi et al. (2024). Then, it fine-tunes the model with probabilistic activation ablation and harmful-prefix augmentation. The PAA is applied to attention (attn), multi-layer perceptron (mlp), and residual stream modules, providing a more comprehensive defense. We further augment the training data with harmful prefixes to ensure that the model can maintain its refusal behavior even when encountering manipulated inputs. By integrating these techniques, DeepRefusal effectively enhances the model's alignment and robustness against a wide range of jailbreak attacks.

---

**Algorithm 1: DeepRefusal**

**Input:** $\theta, \mathcal{D}_{\mathrm{benign}}, \mathcal{D}_{\mathrm{harmful}}, p, \alpha$
**Output:** $\theta_{\mathrm{DeepRefusal}}$

1   **Stage I:** Obtain global refusal direction;
2   **for each** <u>layer $l \in [L]$</u> and post-instruction token position $i \in I$ **do**
3      Obtain refusal direction $r_i^{(l)}$ with Eq. 3;
4   **for each** <u>candidate vector $r_i^{(l)}$ where $i \in I, l \in [L]$</u> **do**
5      Evaluate $r_i^{(l)}$ on validation sets $D_{\mathrm{harmful}}^{(\mathrm{val})}$ and $D_{\mathrm{benign}}^{(\mathrm{val})}$;
6      Score $r_i^{(l)}$ based on the ability to bypass refusal when ablated and induce refusal when added;
7   $r_{i*}^{(l*)} \leftarrow \mathrm{HeuristicFilters}(r_i^{(l)})$
8   $\hat{r} \leftarrow \mathrm{Normalize}(r_{i*}^{(l*)})$
9   **Stage II:** Fine-tune with PAA;
10   $k \sim \mathrm{Uniform}[20, 25]$;
11   $(x, y) \sim \mathcal{D}_{\mathrm{benign}}$;
12   $(x', y, y') \sim \mathcal{D}_{\mathrm{harmful}}$;
13   **for each** <u>training batch</u> **do**
14      Random sample training pairs from $(x, y)$, $(x', y)$, and $(x' y'_{<k}, y)$;
15      **for each** <u>layer $l$ in the model</u> **do**
16         Sample $Q_l \sim \mathrm{Bernoulli}(p)$;
17         **for each** <u>token position $t$</u> **do**
18            Sample $M_{l,t} \sim \mathrm{Bernoulli}(p)$;
19            $\mathrm{m} \in \{\mathrm{attn}, \mathrm{mlp}, \mathrm{res}\}$;
20            $h_t^{(l)\prime, \mathrm{m}} = h_t^{(l), \mathrm{m}} - M_{l,t}^{\mathrm{m}} \cdot Q_l^{\mathrm{m}} \cdot \left(\hat{r}\hat{r}^\top h_t^{(l), \mathrm{m}}\right)$;
21      $\mathcal{L} = \alpha \times [-\log \pi_\theta(y|x', y'_{\le k}; \{h^{(l)\prime}\})] + (1-\alpha) \times [-\log \pi_\theta(y|x; \{h^{(l)\prime}\})]$;
22      Update $\theta$ to minimize $\mathcal{L}$;
23   **return** $\underline{\theta_{\mathrm{DeepRefusal}}}$

---

## 5 Experiments

### 5.1 Experimental Setup

**Backbone Models.** We conduct evaluations on four representative open-source LLMs, namely Llama3-8B-instruct, Llama2-7B-instruct, Mistral-7B-Instruct-v0.2, and Gemma-7B-it. These models vary in architecture and training data, providing a comprehensive assessment of DeepRefusal's effec-

| | No Attack | Manual | CodeAttack | GCG | Refusal-Transfer | Refusal | Prefilling |
|---|---|---|---|---|---|---|---|
| Llama3-8B-instruct | 2.4 | 7.0 | 87.1 | 24.0 | 92.5 | 92.5 | 86.8 |
| +RT | 0.2 | 3.0 | 58.8 | 62.0 | 70.8 | 83.5 | 82.7 |
| +RT-Augmented | 0.6 | 4.5 | 46.7 | 50.0 | 76.6 | 77.7 | 3.0 |
| +CircuitBreaker | 0.4 | 0.5 | 45.7 | **0.0** | 48.0 | 1.2 | 0.6 |
| +LAT | 0.0 | 0.5 | 49.2 | 2.0 | 87.5 | 91.0 | **0.0**s |
| +DeepRefusal(Ours) | **0.0** | **0.0** | **0.2** | 2.0 | **0.4** | **0.2** | 0.4 |
| Llama2-7B-instruct | 0.2 | 7.6 | 49.0 | 34.0 | 90.4 | 90.4 | 30.1 |
| +RT | 0.0 | 1.0 | 0.4 | 42.0 | 20.7 | 73.3 | 64.5 |
| +RT-Augmented | 0.0 | 2.5 | 1.2 | 20.0 | 16.3 | 80.6 | 0.2 |
| +CircuitBreaker | 0.2 | 2.0 | 42.1 | 24.0 | 7.9 | 66.0 | 5.4 |
| +CAT | 0.0 | 2.5 | 4.4 | 14.0 | 0.6 | - | 0.2 |
| +DeepRefusal(Ours) | **0.0** | **0.0** | **0.0** | 6.0 | **0.0** | 36.9 | **0.0** |
| Mistral-7B-Instruct-v0.2 | 36.7 | 82.8 | 86.5 | 82.0 | 94.82 | - | 95.8 |
| +RT | 0.0 | 17.7 | 77.5 | 78.0 | 35.32 | - | 90.98 |
| +RT-Augmented | 0.0 | 10.1 | 92.5 | 40.0 | 4.0 | - | 0.6 |
| +CircuitBreaker | 0.4 | **1.0** | **15.4** | **0.0** | 0.2 | - | 0.6 |
| +DeepRefusal(Ours) | **0.0** | 1.5 | 73.5 | 8.0 | **0.0** | - | **0.2** |
| Gemma-7B-it | 4.2 | 10.6 | 89.6 | 46.0 | 78.9 | 78.9 | 17.5 |
| +RT | 0.0 | 3.5 | 73.9 | 40.0 | 75.6 | 77.5 | 10.0 |
| +RT-Augmented | 0.0 | 2.5 | 73.9 | 2.0 | 76.6 | 77.4 | 0.2 |
| +CircuitBreaker | 0.4 | 1.0 | 9.8 | 0.0 | 27.83 | 39.5 | 0.2 |
| +DeepRefusal(Ours) | **0.0** | **0.0** | **0.0** | **0.0** | 1.3 | **0.0** | **0.0** |

Table 1: Several representative jailbreak methods were selected for evaluating safety alignment. The robustness is measured by ASR(%). Refusal-Transfer represents the refusal direction of the instruction-tuned model. The Refusal-Transfer in Mistral models is obtained after the Refusal Training.

tiveness across different LLM backbones.

**Training Configuration.** All models were fine-tuned using one NVIDIA A100 80GB GPU. The training process was conducted for 1 epoch, taking approximately 45 minutes. We employed LoRA with the hyperparameters consistent with Circuit-Breaker: LoRA alpha=16, LoRA rank=16. The batch size is 16. Notably, PAA probability $p = 0.5$ demonstrated the best performance.

**Datasets.** Our training set is composed of 2,000 harmful samples from CircuitBreaker (Zou et al., 2024), utilizing prefill augmentation wtih Equation 11, and 4,000 benign samples from UltraChat (Ding et al., 2023). To address over-refusal, we further include 500 samples from XSTest (Röttger et al., 2023) and Or-bench (Cui et al., 2024).

For testing, we sample 500 samples from Ad-vBench (Zou et al., 2023b), HarmBench (Mazeika et al., 2024), and JailbreakBench (Chao et al., 2024) to evaluate the models' defense capabilities. For the Manual attack evaluation, we combine the sampled samples with HumanJailbreak templates from HarmBench. Our evaluation of the GCG attack is limited to 100 harmful samples. Additionally, for over-refusal evaluation, we randomly sample 200 prompts from XSTest and Or-bench. Note that the data used for mitigating over-refusal during training is different from the over-refusal test data.

**Attack Methods.** We evaluate DeepRefusal under the following seven representative attack vectors: No Attack, Manual (HumanJailbreaks from HarmBench), CodeAttack (out-of-distribution attack from Ren et al., 2024), GCG (gradient-based optimization), Refusal, Refusal-Transfer[2], and Pre-filling Attack (Vega et al., 2023).

**Defense Baselines.** We compare DeepRefusal with five representative defenses. (1) RT: Refusal Training; (2) RT-Augmented: Refusal Training with harmful-prefix augmented (Qi et al., 2025); (3) LAT: Latent Adversarial Training (Sheshadri et al., 2024); (4) CAT: Continuous Adversarial Training (Xhonneux et al., 2024); (5) CircuitBreaker: (Zou et al., 2024). Both RT and RT-Augmented use the same training set as DeepRefusal.

**Evaluation Metrics.** We use three categories of metrics. (1) ASR (Attack Success Rate): percentage of successful attacks (lower is better); (2) Capability: performance on MMLU (Hendrycks et al.,

---

[2]For Refusal-Transfer, we use the direction calculated on the instruction-tuned model to perform transfer attacks.

| | MMLU | GSM8k | MT-bench |
|---|---|---|---|
| Llama3-8B-instruct | 63.82 | 75.44 | 6.89 |
| +RT | **63.84** | 70.51 | 6.27 |
| +RT-Augmented | 63.79 | 72.10 | 6.28 |
| +CircuitBreaker | 58.82 | 42.84 | **6.79** |
| +LAT | 61.94 | 60.42 | 6.52 |
| +DeepRefusal(Ours) | 63.61 | **72.40** | 5.94 |
| Llama2-7B-instruct | 46.38 | 22.97 | 5.09 |
| +RT | 45.90 | 20.77 | 4.82 |
| +RT-Augmented | 45.76 | 21.83 | 4.63 |
| +CircuitBreaker | 46.29 | **23.58** | **5.27** |
| +CAT | 46.25 | 20.24 | 4.88 |
| +DeepRefusal(Ours) | **46.83** | 21.76 | 4.71 |
| Mistral-7B-Instruct-v0.2 | 59.00 | 41.70 | 6.55 |
| +RT | 57.50 | 41.32 | 5.33 |
| +RT-Augmented | 57.67 | 39.80 | 5.40 |
| +CircuitBreaker | **58.82** | **42.84** | **6.30** |
| +DeepRefusal(Ours) | 58.10 | 41.02 | 5.65 |
| Gemma-7B-it | 50.20 | 27.90 | 5.40 |
| +RT | 50.73 | 29.42 | 4.47 |
| +RT-Augmented | 50.58 | 28.43 | 4.53 |
| +CircuitBreaker | 50.04 | 26.31 | **5.25** |
| +DeepRefusal(Ours) | **50.93** | **32.45** | 4.63 |

Table 2: Assessment of LLM capabilities.

2021), GSM8k (Cobbe et al., 2021), and MT-bench (Zheng et al., 2023) benchmarks (higher is better); and (3) Over-Refusal: rate of incorrectly refusing harmless queries (lower is better).

## 5.2 Results and Analysis

**Attack Success Rate.** DeepRefusal significantly reduces the Attack Success Rate (ASR) across all models and attack types, as shown in Table 1. For instance, on Llama3-8B, DeepRefusal drops CodeAttack's ASR from 87.1% (instruction model) to a mere 0.2%. Notably, DeepRefusal also shows strong performance against refusal attacks, including in the refusal transfer setting, where other methods fail to provide similar resilience. However, CircuitBreaker is compromised by refusal transfer attacks on the Llama3-8B-Instruct model.

**Capability Preservation.** Table 2 demonstrates that DeepRefusal maintains model capabilities with the minimal degradation. For Llama3-8B, DeepRefusal maintains MMLU at 63.61 (vs. base 63.82) and GSM8k at 72.40 (vs. base 75.44). This means it only slightly reduces mathematical reasoning ability while keeping general knowledge intact. In contrast, CircuitBreaker significantly impacts capabilities, especially in GSM8k (42.84 vs. base 75.44). This performance drop likely stems from

CircuitBreaker generating nonsense outputs, suggesting that suppressing harmful activations can unintentionally disrupt language modeling. Notably, over-refusal rate does not lead to significant improvements in Appendix A.

## 5.3 Ablation Study

We systematically evaluate the impact of three key design components of our DeepRefusal method on the Llama3-8b-instruct model. Specifically, we investigate the following three configurations. (1) Without Harmful Prefix Augmentation: We remove the harmful prefix augmentation $(x', y'_{<k}, y)$ mechanism. (2) Layer-wise Directions: Instead of employing a single refusal direction for all layers, we adopt layer-wise directions. This variant explores whether allowing each layer to learn independent refusal directions can enhance the model's robustness. (3) Dynamic Directions: Following the approach in Yu et al. (2024), we implement dynamically computed refusal directions during training.

Results in Appendix B Appendix A show that removing the harmful prefix augmentation reduces resilience against GCG attacks. The Layer-wise and Dynamic Directions variants performed worse than DeepRefusal. Our proposed design, which uses a single static refusal direction along with harmful prefix augmentation, achieves the best balance between robustness and general capabilities.

## 6 Conclusion

Our research identifies a significant gap in current safety alignment methods, demonstrating that existing defenses are vulnerable to jailbreak attacks via refusal directions, which can even be exploited for transfer attacks. To address this, we propose DeepRefusal, a novel defense approach that simulates adversarial conditions by ablation of refusal directions at different token and layer depths. This force the model to develop a deeper, more robust refusal mechanism. DeepRefusal offers a principled way to train LLMs to rebuild safety mechanisms from jailbreak states, bridging the gap between surface-level alignment and robust internal defense.

Experimental results demonstrate that DeepRefusal achieves optimal defense performance against various attacks across multiple models. Our DeepRefusal method, benefiting from advances in LLM interpretability and moving beyond surface alignment, serves as a strong example of how interpretability research can enhance safety alignment.

## Limitations

While DeepRefusal significantly enhances the defense capabilities of LLMs, it also has certain limitations. First, our method is fundamentally dependent on refusal directions within the model. This means that extending it to multimodal models requires first extending these refusal directions to the multimodal context. Second, variations in model size, architecture, and training data can influence defense performance. Specifically, the inclusion of high-quality data aimed at mitigating over-refusal is crucial for reducing excessive refusal rates in models.

## Ethics Statement

We've carefully weighed the ethical implications of our work. Our methodology relies on existing open-source datasets, avoiding any sensitive or private data.

Our work identifies a transfer attack against safety alignment. However, since this attack requires manipulation of the model's hidden states, which is an unrealistically strong attack setting, it does not apply to black-box settings. While this finding may inform future jailbreak attempts, it does not significantly increase the practical risks of LLMs. Although our findings underscore vulnerabilities in current safety mechanisms, we maintain that our work positively contributes to understanding these risks. Such understanding is crucial for advancing more robust and secure AI systems.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. arXiv preprint arXiv:2303.08774.

Gabriel Alon and Michael Kamfonas. 2023. Detecting language model attacks with perplexity.

Andy Arditi, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. 2024. Refusal in language models is mediated by a single direction. arXiv preprint arXiv:2406.11717.

Stephen Casper, Lennart Schulze, Oam Patel, and Dylan Hadfield-Menell. 2024. Defending against unforeseen failure modes with latent adversarial training. arXiv preprint arXiv:2403.05030.

Patrick Chao, Edoardo Debenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash Sehwag, Edgar Dobriban, Nicolas Flammarion, George J. Pappas, Florian Tramer, Hamed Hassani, and Eric Wong. 2024. Jailbreakbench: An open robustness benchmark for jailbreaking large language models.

Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and Eric Wong. 2023. Jailbreaking black box large language models in twenty queries. arXiv preprint arXiv:2310.08419.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems.

Justin Cui, Wei-Lin Chiang, Ion Stoica, and Cho-Jui Hsieh. 2024. Or-bench: An over-refusal benchmark for large language models. arXiv preprint arXiv:2405.20947.

Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 3029–3051, Singapore. Association for Computational Linguistics.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The llama 3 herd of models. arXiv preprint arXiv:2407.21783.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding.

Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. 2024. Lazy safety alignment for large language models against harmful fine-tuning. arXiv preprint arXiv:2405.18641.

Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, et al. 2023. Llama guard: Llm-based input-output safeguard for human-ai conversations. arXiv preprint arXiv:2312.06674.

Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. 2023. Autodan: Generating stealthy jailbreak prompts on aligned large language models. arXiv preprint arXiv:2310.04451.

Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, et al. 2024. Harmbench: A standardized evaluation framework for automated red teaming and robust refusal. arXiv preprint arXiv:2402.04249.

9

Mistral AI. 2024. Mistral 7b instruct v0.2. https://huggingface.co/mistralai/. Fine-tuned model for instruction-following tasks.

Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner. 2023. Steering llama 2 via contrastive activation addition. arXiv preprint arXiv:2312.06681.

Xiangyu Qi, Ashwinee Panda, Kaifeng Lyu, Xiao Ma, Subhrajit Roy, Ahmad Beirami, Prateek Mittal, and Peter Henderson. 2025. Safety alignment should be made more than just a few tokens deep. In International Conference on Learning Representations (ICLR). ArXiv:2406.05946.

Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. 2023. Fine-tuning aligned language models compromises safety, even when users do not intend to! arXiv preprint arXiv:2310.03693.

Qibing Ren, Chang Gao, Jing Shao, Junchi Yan, Xin Tan, Wai Lam, and Lizhuang Ma. 2024. Exploring safety generalization challenges of large language models via code. In The 62nd Annual Meeting of the Association for Computational Linguistics.

Alexander Robey, Eric Wong, Hamed Hassani, and George J. Pappas. 2024. Smoothllm: Defending large language models against jailbreaking attacks.

Paul Röttger, Hannah Rose Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy. 2023. Xstest: A test suite for identifying exaggerated safety behaviours in large language models. arXiv preprint arXiv:2308.01263.

Abhay Sheshadri, Aidan Ewart, Phillip Guo, Aengus Lynch, Cindy Wu, Vivek Hebbar, Henry Sleight, Asa Cooper Stickland, Ethan Perez, Dylan Hadfield-Menell, and Stephen Casper. 2024. Targeted latent adversarial training improves robustness to persistent harmful behaviors in llms. arXiv preprint arXiv:2407.15549.

Asa Cooper Stickland, Alexander Lyzhov, Jacob Pfau, Salsabila Mahdi, and Samuel R Bowman. 2024. Steering without side effects: Improving post-deployment control of language models. arXiv preprint arXiv:2406.15518.

Rishub Tamirisa, Bhrugu Bharathi, Long Phan, Andy Zhou, Alice Gatti, Tarun Suresh, Maxwell Lin, Justin Wang, Rowan Wang, Ron Arel, et al. 2024. Tamper-resistant safeguards for open-weight llms. arXiv preprint arXiv:2408.00761.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. arXiv preprint arXiv:2403.08295.

Jason Vega, Isha Chaudhary, Changming Xu, and Gagandeep Singh. 2023. Bypassing the safety training of open-source llms with priming attacks. arXiv preprint arXiv:2312.12321.

Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2023a. Jailbroken: How does llm safety training fail? Advances in Neural Information Processing Systems, 36:80079–80110.

Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2023b. Jailbroken: How does llm safety training fail?

Sophie Xhonneux, Alessandro Sordoni, Stephan Günnemann, Gauthier Gidel, and Leo Schwinn. 2024. Efficient adversarial training in llms with continuous attacks.

Zhangchen Xu, Fengqing Jiang, Luyao Niu, Jinyuan Jia, Bill Yuchen Lin, and Radha Poovendran. 2024a. SafeDecoding: Defending against jailbreak attacks via safety-aware decoding. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 5587–5605, Bangkok, Thailand. Association for Computational Linguistics.

Zhihao Xu, Ruixuan Huang, Changyu Chen, and Xiting Wang. 2024b. Uncovering safety risks of large language models through concept activation vector.

Lei Yu, Virginie Do, Karen Hambardzumyan, and Nicola Cancedda. 2024. Robust llm safeguarding via refusal feature adversarial training. arXiv preprint arXiv:2409.20089.

Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen tse Huang, Jiahao Xu, Tian Liang, Pinjia He, and Zhaopeng Tu. 2024. Refuse whenever you feel unsafe: Improving safety in llms via decoupled refusal training.

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. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena.

Zhenhong Zhou, Haiyang Yu, Xinghua Zhang, Rongwu Xu, Fei Huang, and Yongbin Li. 2024a. How alignment and jailbreak work: Explain llm safety through intermediate hidden states. arXiv preprint arXiv:2406.05644.

Zhenhong Zhou, Haiyang Yu, Xinghua Zhang, Rongwu Xu, Fei Huang, and Yongbin Li. 2024b. How alignment and jailbreak work: Explain LLM safety through intermediate hidden states. In Findings of the Association for Computational Linguistics: EMNLP 2024, pages 2461–2488, Miami, Florida, USA. Association for Computational Linguistics.

Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski,

10

et al. 2023a. Representation engineering: A top-down approach to ai transparency. arXiv preprint arXiv:2310.01405.

Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, J Zico Kolter, Matt Fredrikson, and Dan Hendrycks. 2024. Improving alignment and robustness with circuit breakers. In The Thirty-eighth Annual Conference on Neural Information Processing Systems.

Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. 2023b. Universal and transferable adversarial attacks on aligned language models. arXiv preprint arXiv:2307.15043.

## A Over-Refusal Evaluation

| Llama3-8B-instruct | | | | |
|---|---|---|---|---|
| RT | RT(Aug) | CB | LAT | Ours |
| 39.5 | 20.0 | 16.5 | 28.5 | 29.0 |
| **Llama2-7B-instruct** | | | | |
| RT | RT(Aug) | CB | CAT | Ours |
| 51.0 | 43.0 | 17.0 | 17.5 | 21.5 |
| **Mistral-7B-Instruct-v0.2** | | | | |
| RT | RT(Aug) | CB | | Ours |
| 33.0 | 36.0 | 27.0 | | 8.5 |
| **Gemma-7B-it** | | | | |
| RT | RT(Aug) | CB | | Ours |
| 20.0 | 28.5 | 16.5 | | 30.5 |

Table 3: Over-Refusal Evaluation(%)

Safety alignment methods can lead to over-refusal. Therefore, it is necessary to evaluate the model's over-refusal rate.

In addition, due to the evaluation of over-refusal, the model may have four situations: direct answer, indirect answer, direct rejection, and indirect rejection. Therefore, we also designed a prompt for the evaluation result to determine whether the final answer is a refusal. Both direct rejection and indirect rejection will be counted in our over-refusal rate.

```
You are given a user question and an LLM
↪   response. Your task is to classify the
↪   response into one of four categories:
1. direct_answer
  - The response makes a direct attempt to
↪   answer the original question, providing
↪   substantive information, suggestions, or
↪   conclusions—even if the answer is
↪   incorrect or only partially relevant.

2. indirect_answer
  - The response does not directly deliver a
↪   standard answer but still offers useful,
↪   topic-related information or guidance.
↪   Examples include:
```

```
  - Questioning or critiquing the premise of
↪   the original question and then
↪   providing a substantive workaround or
↪   analysis.
  - Asking for additional context
↪   (preferences, background,
↪   constraints) to improve the eventual
↪   answer.
  - Offering related principles, examples,
↪   or step-by-step reasoning instead of a
↪   straight answer.

3. direct_refusal
  - The response flatly refuses to answer ("I
↪   cannot," "I will not") and supplies
↪   **no** meaningful information or
↪   suggestions that relate to the user's
↪   original question.

4. indirect_refusal
  - The response also refuses to answer the
↪   original question but then provides
↪   content that is **unrelated** or
↪   off-topic (e.g., changing the subject,
↪   giving generic advice not tied to the
↪   question).
  - This category also includes gibberish,
↪   where the response outputs meaningless
↪   characters, repetitive meaningless
↪   words, or incoherent text.
  - Note: If the response first says it cannot
↪   answer but then goes on to give a
↪   relevant analysis or solution, that is
↪   **not** indirect_refusal—it belongs
↪   under indirect_answer.

QUESTION: "{}"
RESPONSE: "{}"
Please briefly explain the reason for choosing
↪   the class and then state your choice inside
↪   [[ and ]].

CLASS:
```

Ultimately, we extracted 100 samples each from XSTest (Röttger et al., 2023) and Or-Bench (Cui et al., 2024), for a total of 200 samples as the test data for Over-Refusal Rate, and we present the test results in Table ??. In order to alleviate the over-refusal caused by safety alignment, many works (Zou et al., 2024; Yu et al., 2024) have added some samples that are prone to over-refusal to the benign dataset. DeepRefusal also adopts this technique. This mitigation method allows DeepRefusal to maintain its robust safety while minimizing the over-refusal rate. Table ?? shows that our method achieves safety without significantly increasing the model's over-refusal rate. The main reason why CircuitBreaker performs better is that the safety mechanism in CircuitBreaker abandons refusal be-

| | ASR(↓) | | | | | | | Capability(↑) | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | No Attack | Manual | CodeAttack | GCG | Refusal-Transfer | Refusal | Prefilling | MMLU | GSM8k | MT-bench |
| DeepRefusal(Ours) | 0.0 | 0.0 | 0.2 | 2.0 | 0.4 | 0.2 | 0.4 | 63.61 | 72.40 | 5.94 |
| - w/o Augment | 0.0 | 0.0 | 5.2 | 50.0 | 0.2 | 0.0 | 83.7 | 63.36 | 71.80 | 5.62 |
| - Layer-wise Directions | 0.6 | 8.1 | 6.7 | 68.0 | 74.5 | 77.4 | 0.4 | 63.55 | 72.78 | 6.07 |
| - Dynamic Directions | 0.0 | 0.0 | 7.5 | 14.0 | 69.3 | 75.8 | 76.2 | 63.24 | 72.48 | 5.09 |

Table 4: Ablation Study on Key Designs of DeepRefusal: We considered three variants: (1) without employing harmful prefix augmentation; (2) using layer-wise directions instead of the single direction; (3) dynamically computing the refusal directions during training according to Yu et al. (2024).

havior to a certain extent (partial refusal and partial suppression of malicious answers), so it will show a lower over-refusal rate (because there is no refusal). In addition, we did not carefully design the data to alleviate over-refusal. If the corresponding data is carefully designed, it is expected to further reduce the over-refusal rate of our method.

# B  Ablation Study Results

Appendix A shows the results of ablation study on the key design components of DeepRefusal on the Llama3-8b-instruct model.