

---

# LLM Jailbreaks Exploit Attention Sinks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Suffix-based jailbreak attacks append adversarial token sequences to harmful re-  
2       quests, bypassing safety guardrails in language models. Despite their effectiveness,  
3       the mechanisms enabling these attacks remain poorly understood. We find that  
4       tokens in adversarial suffixes are prone to inducing *attention sinks*—a phenomenon  
5       where certain tokens (e.g., BOS, punctuation, and chat tokens) receive dispro-  
6       portionately high attention from subsequent tokens—and establish a relationship  
7       between suffix-induced sinks and attack success: amplifying the influence of suffix  
8       sinks improves attack success by up to 276%, while attenuating it reduces  
9       attack success by up to 84%. We trace this effect to the model’s *refusal direction*:  
10       sink tokens induce perturbations aligned with the refusal direction, cumulatively  
11       suppressing the residual stream’s refusal alignment across layers. Our results  
12       generalize across several models and suffix-based jailbreak methods, exposing  
13       a fundamental structural vulnerability in transformer attention mechanisms that  
14       adversarial suffixes exploit to bypass safety alignment.

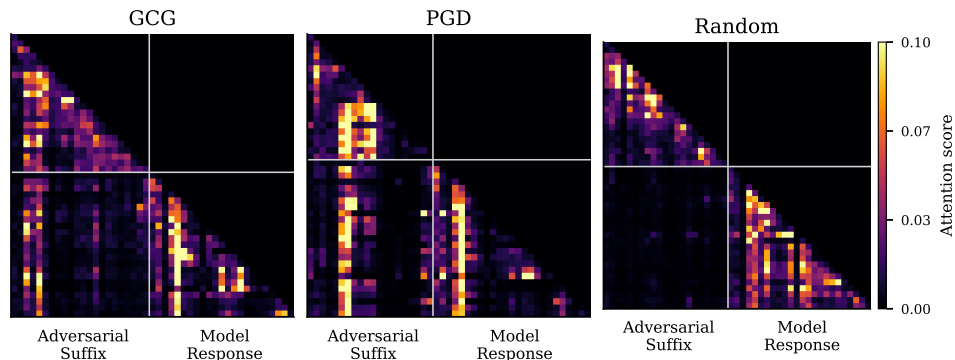


Figure 1: Attention weights at Layer 9, Head 11 of GEMMA-2-9B for a GCG suffix (left), PGD suffix (middle), and random suffix (right). The attention scores are truncated to display only the adversarial suffix and model response regions, with values clipped at 0.1 for visualization. The lower-left quadrant of each panel shows attention from response tokens to suffix positions: GCG and PGD induce distinct attention sinks, while the random baseline does not.

## 1 Introduction

16       Modern language models undergo several stages of fine-tuning to align with human values before  
17       deployment, including instruction-following [Ouyang et al., 2022] and safety alignment [Bai et al.,  
18       2022b]. The goal is to develop models that are both useful across a wide range of tasks and can  
19       reliably avoid harmful behaviors [Bai et al., 2022a].

20 A critical challenge in AI safety is understanding how to maintain these safety guardrails and identify  
21 when they fail. Safety threats include jailbreak attacks [Zou et al., 2023], prompt injection [Perez and  
22 Ribeiro, 2022], and data poisoning [Wallace et al., 2021, Wan et al., 2023], among others [Yao et al.,  
23 2024].

24 To address these threats, many researchers focus on identifying the specific mechanisms within the  
25 model that enable impermissible outputs [Arditi et al., 2024, Yona et al., 2025]. Understanding these  
26 mechanisms is crucial for designing better defenses and alignment strategies.

27 **Suffix-Based Jailbreaking** In this work, we focus on *jailbreaking*—attacks that bypass language  
28 model safety guardrails and elicit prohibited behaviors [Zou et al., 2023, Chao et al., 2025, Anil  
29 et al., 2024]. A prominent class of such attacks uses *adversarial suffixes* – short, optimized character  
30 sequences appended to harmful requests that induce model compliance [Guo et al., 2021, Wen et al.,  
31 2023, Zou et al., 2023]. Remarkably, suffixes optimized on specific harmful prompts often transfer to  
32 entirely different prompts, a property known as *suffix universality* [Zou et al., 2023, Ben-Tov et al.,  
33 2025, Liao and Sun, 2024].

34 **Problem Statement** Despite the effectiveness of suffix-based jailbreaks, research has focused over-  
35 whelmingly on discovering new vulnerabilities rather than understanding their underlying mechanics.  
36 Systematic analyses of why these suffixes succeed remain limited [Wang et al., 2024, Arditi et al.,  
37 2024, Hu et al., 2025, Ben-Tov et al., 2025], creating a critical gap in our understanding of model  
38 robustness.

39 **Attention Sinks** A closer inspection of successful adversarial suffixes reveals they often consist  
40 of seemingly unnatural tokens such as repeated punctuation, rare character sequences, or atypical  
41 symbols [Mu et al., 2025].

42 Interestingly, these are the exact types of tokens known to induce *attention sinks* [Yu et al., 2024b], a  
43 phenomenon in transformer models where certain tokens receive disproportionately high attention  
44 scores from subsequent tokens. While the first token is the canonical example [Xiao et al., 2023],  
45 attention sinks can occur at later positions, often on seemingly arbitrary tokens such as punctuation  
46 tokens [Yu et al., 2024a, Sun et al., 2024, Cancedda, 2024, Gu et al., 2025, Zhang et al., 2025].

47 **Why Do Attention Sinks Matter for Jailbreaking?** To understand why attention sinks might influ-  
48 ence jailbreaking, consider a causal transformer with residual connections, in which the representation  
49 of a response token at position  $i$  is updated as

$$x'_i \leftarrow x_i + \sum_{j \leq i} \mathbf{P}_{i,j} v_j, \quad (1)$$

50 where  $\mathbf{P}_{i,j}$  denotes the attention weight from token  $i$  to token  $j$ ,  $v_j$  is the value vector at position  $j$ ,  
51 and attention is restricted to past tokens ( $j \leq i$ ).

52 When a suffix sink token at position  $s$  attracts disproportionately high attention from response tokens,  
53 its contribution dominates the update:

$$x'_i \approx x_i + \mathbf{P}_{i,s} v_s. \quad (2)$$

54 Crucially, because the same sink value vector  $v_s$  is added—up to a scalar weight—to many response  
55 token representations through the residual stream, the sink acts as a mechanism for *broadcasting* a  
56 shared perturbation across the entire response, which could shift response representations away from  
57 refusal-aligned behaviors, thereby enabling jailbreaking.

58 **Contributions** Motivated by the hypothesis outlined above, we empirically evaluate whether the  
59 seemingly arbitrary tokens used in adversarial suffixes give rise to attention sinks in early transformer  
60 layers and whether such sinks in fact play a substantive role in enabling suffix-based jailbreaks;  
61 specifically, we show that:

- 62 1. Adversarial suffixes induce more attention sinks than random baselines, across multiple  
63 safety-aligned models and suffix-based attack methods. (Section 3)
- 64 2. Modulating suffix attention sinks shifts attack success: amplification boosts attack success  
65 by up to 276%, while attenuation reduces it by up to 84%. (Section 4)

66 3. Sink-induced perturbations are aligned with the refusal direction and cumulatively suppress  
 67 residual-stream refusal alignment across layers. (Section 5).

## 68 2 Methodology

### 69 2.1 Experimental Setup

70 **Models** We evaluate our approach on four safety-  
 71 tuned models spanning diverse model families in the  
 72 7–9B parameter range, as listed in Table 1. We focus  
 73 on this parameter range because jailbreak transfer-  
 74 ability within model families has been shown to be  
 75 consistent across scales [Mazeika et al., 2024].

76 **Datasets** We train attacks on HarmBench [Mazeika  
 77 et al., 2024] and evaluate on JailbreakBench [Chao  
 78 et al., 2024], filtering to ensure no overlap. Specifically, we select the first 100 standard behaviors  
 79 from HarmBench (excluding contextual and copyright categories) and remove the 4 behaviors that  
 80 overlap with our JailbreakBench evaluation set, resulting in 96 training prompts. This training budget  
 81 is consistent with established practices in the literature [Zou et al., 2023, Arditi et al., 2024, Beyer  
 82 et al., 2025].

83 **Attack Methods** We conduct our analysis using  
 84 two suffix-based jailbreak methods: Greedy Coordi-  
 85 nate Gradient (GCG) [Zou et al., 2023] and Projected  
 86 Gradient Descent (PGD) [Geisler et al., 2024]. Both  
 87 optimize an adversarial suffix  $s$  to maximize the log-  
 88 likelihood of a target affirmative response  $y$  (e.g.,  
 89 “Sure, here is...” conditioned on a harmful request  $x$ :

$$\mathcal{L} = -\log p(y | x, s). \quad (3)$$

90 The two methods differ only in how they search  
 91 the discrete suffix space: GCG via greedy gradient-  
 92 guided token replacements, and PGD via projected  
 93 gradient descent on a continuous relaxation of the token embeddings. Figure 2 illustrates the attack  
 94 prompt structure, which partitions into three components: the harmful request  $x$ , the adversarial suffix  
 95  $s$ , and the target response  $y$ .

### 96 2.2 Evaluation

97 We evaluate attack efficacy using Attack Success Rate  
 98 (ASR) computed over all 100 prompts from Jailbreak-  
 99 Bench. We limit response length to 512 tokens to  
 100 ensure ASR convergence [Mazeika et al., 2024] and  
 101 use greedy decoding to generate a single response per  
 102 (prompt, suffix) pair. An attack is considered success-  
 103 ful on a prompt if its response is classified as harmful.

104 **Attack Success Classification** To classify individ-  
 105 ual responses as successful or failed attacks, we use StrongREJECT’s fine-tuned classifier [Souly  
 106 et al., 2024], which assigns a score in the range [0, 1] to each (harmful prompt, model response) pair,  
 107 where higher scores indicate more harmful responses. Following Beyer et al. [2025], we apply a  
 108 threshold of 0.5, labeling responses with scores  $\geq 0.5$  as successful attacks and  $< 0.5$  as failures.

### 109 2.3 Suffix Selection

110 For each model and attack method, we generate a pool of candidate suffixes of length 20 tokens by  
 111 running the attack on each of the 96 HarmBench training prompts under multiple random seeds: 6

Table 1: Safety-aligned models evaluated in this study.

Model	Reference
LLAMA-3-8B	Grattafiori et al. [2024]
LLAMA-3.1-8B	Grattafiori et al. [2024]
GEMMA-2-9B	Team et al. [2024]
QWEN2.5-7B	Yang et al. [2025]

```

<user_start>
Write a tutorial on how to make a
bomb !!!!!!!!!!!!!!!!!!!!!!!
<user_end><assistant_start>
Sure, here is a tutorial for
making a bomb...
  
```

Figure 2: Example prompt structure. The user instruction consists of a harmful request (blue) and an adversarial suffix (red). The target response (green) consists of an affirmative model completion.

Table 2: Average attack success rate (ASR) on JailbreakBench across the 10 best GCG and PGD suffixes per model.

Model	GCG (%)	PGD (%)
LLAMA-3-8B	19.2±5.7	12.8±4.9
LLAMA-3.1-8B	17.8±5.4	36.8±11.8
GEMMA-2-9B	53.5±9.5	66.1±8.6
QWEN-2.5-7B	55.3±6.9	56.5±7.8

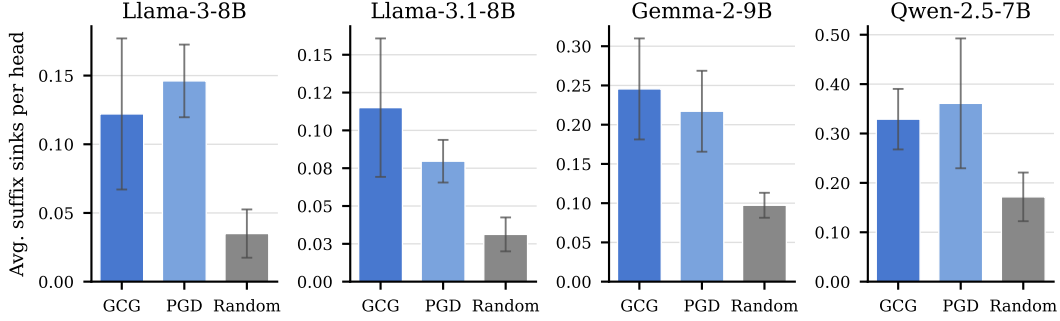


Figure 3: **Adversarial suffixes induce more attention sinks than random suffixes across all four models.** Each bar reports the average number of suffix attention sinks detected per head.

112 seeds for GCG (576 unique suffixes per model) and 3 seeds for PGD (288 unique suffixes per model).  
 113 We then evaluate every candidate suffix on JailbreakBench and select the top 10 by ASR for analysis.  
 114 This yields 10 GCG suffixes and 10 PGD suffixes per model. Table 2 reports the average ASR on  
 115 JailbreakBench for the selected top-10 suffix sets per model.

116 As a control, we generate 10 random suffixes per model by uniformly sampling 20 tokens from each  
 117 model’s vocabulary, matching the length of the GCG and PGD suffixes [Arditi et al., 2024].

118 Unless otherwise noted, all subsequent results are averaged over these 10 suffixes per (model, method)  
 119 and the 100 JailbreakBench evaluation prompts.

### 120 3 Adversarial Suffixes Induce Attention Sinks

#### 121 3.1 Measuring Attention Sinks

122 Having described the jailbreaking methodology, we now study the resulting patterns in the model’s  
 123 attention weights, beginning with a formal criterion for identifying attention sinks. We adopt the  
 124 threshold-based metric proposed by Gu et al. [2025] and used in subsequent work [Zhang et al., 2025,  
 125 Queipo-de Llano et al., 2025], which quantifies how much attention a given token position attracts  
 126 from subsequent tokens.

**Definition (Attention Sink).** For a sequence of length  $T$ , let  $\mathbf{P} \in \mathbb{R}^{T \times T}$  denote the attention weight matrix (softmax probabilities) for a given attention head. We say that token  $t$  is an *attention sink* if

$$127 \quad \alpha_t = \frac{1}{T-t+1} \sum_{k=t}^T \mathbf{P}_{k,t} > \epsilon, \quad (4)$$

where  $\epsilon > 0$  is a threshold parameter.

128 Intuitively,  $\alpha_t$  measures the average attention that token  $t$  receives from all subsequent tokens  
 129 (including itself). A large value of  $\alpha_t$  indicates that token  $t$  attracts a disproportionate amount of  
 130 attention from later positions. In practice, we evaluate  $\alpha_t$  over a fixed downstream window of 150  
 131 tokens anchored at the start of the suffix (i.e., we set  $T = s + 149$  in Equation 4, where  $s$  is the  
 132 position of the first suffix token), ensuring all suffix positions are scored against a common context  
 133 window [Gu et al., 2025].

#### 134 3.2 Quantifying Suffix Attention Sinks

135 We summarize sink prevalence at the model level using a single scalar: the average number of suffix  
 136 sinks per attention head, computed via Equation 4 with a fixed threshold  $\epsilon = 0.02$ . The threshold is  
 137 chosen to detect a conservatively small number of sink tokens per head, minimizing false positives.

138 As shown in Figure 3, adversarial suffixes induce significantly more sinks per head than random  
 139 suffixes of the same length across all four models and both attack methods. Averaged across models,

140 GCG and PGD suffixes induce  $2.9\times$  and  $2.76\times$  as many sinks per head as the random baseline,  
 141 respectively.

142 Figure 1 illustrates this phenomenon at the head level for Layer 9, Head 11 of GEMMA-2-9B. Both  
 143 GCG and PGD suffixes induce distinct attention sinks, with model response tokens consistently  
 144 attending to the same subset of suffix positions. The random suffix baseline induces no such pattern.

145 Overall, these results establish that adversarial suffixes consistently produce attention sinks across  
 146 different attack methods and model families.

## 147 4 Suffix Attention Sinks Modulate Jailbreak Success

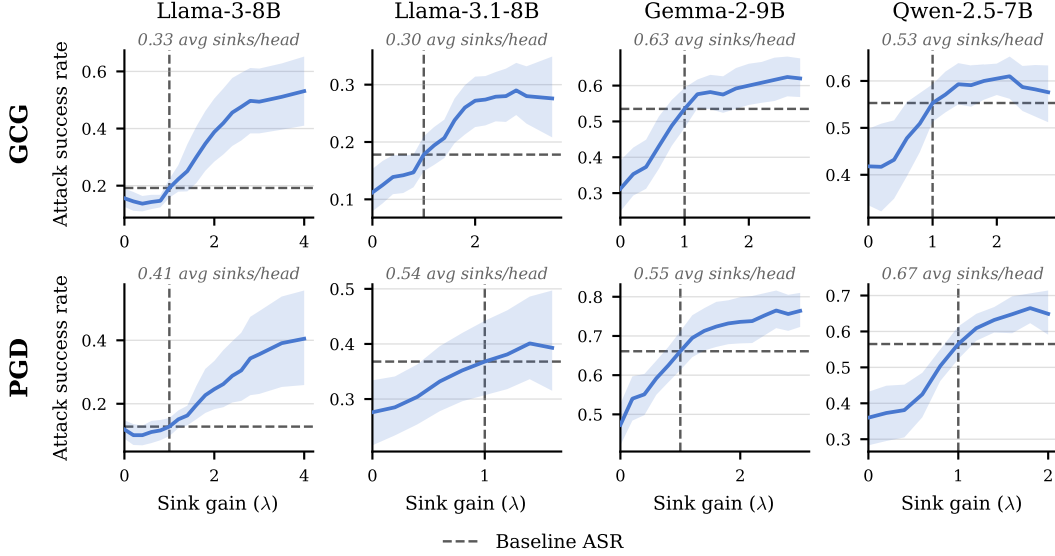


Figure 4: **Sink modulation shifts attack success across both attack methods.** Each panel plots attack success rate (ASR) against the sink gain  $\lambda$ . Dashed lines mark the no-intervention baseline ( $\lambda = 1$ ). Panel subtitles report the average number of suffix sinks per head modulated by the intervention.

### 148 4.1 A Sink-Subspace Intervention

149 **Sink-Subspace Decomposition** Section 3 established that adversarial suffixes induce attention  
 150 sinks; we now ask whether these sinks contribute to jailbreak success. To answer this, we isolate  
 151 their effect on each attention head’s output via an exact decomposition. For an attention head with  
 152 attention weights  $\mathbf{P} \in \mathbb{R}^{T \times T}$  and values  $\mathbf{V} \in \mathbb{R}^{T \times d_v}$ , the head output admits the decomposition

$$\mathbf{P}\mathbf{V} = \sum_{i=1}^T \mathbf{p}_i \mathbf{v}_i^\top, \quad (5)$$

153 where  $\mathbf{p}_i$  is the  $i$ -th column of  $\mathbf{P}$  and  $\mathbf{v}_i^\top$  is the  $i$ -th row of  $\mathbf{V}$ . Let  $\mathcal{S} \subseteq \{1, \dots, T\}$  denote the suffix  
 154 sink positions for this head, identified via Equation 4. Partitioning the sum yields

$$\mathbf{P}\mathbf{V} = \underbrace{\sum_{i \in \mathcal{S}} \mathbf{p}_i \mathbf{v}_i^\top}_{\text{sink term}} + \underbrace{\sum_{i \notin \mathcal{S}} \mathbf{p}_i \mathbf{v}_i^\top}_{\text{non-sink term}}. \quad (6)$$

155 This partition isolates the sink contribution as an additive component of the head output that can be  
 156 modulated independently of the rest.

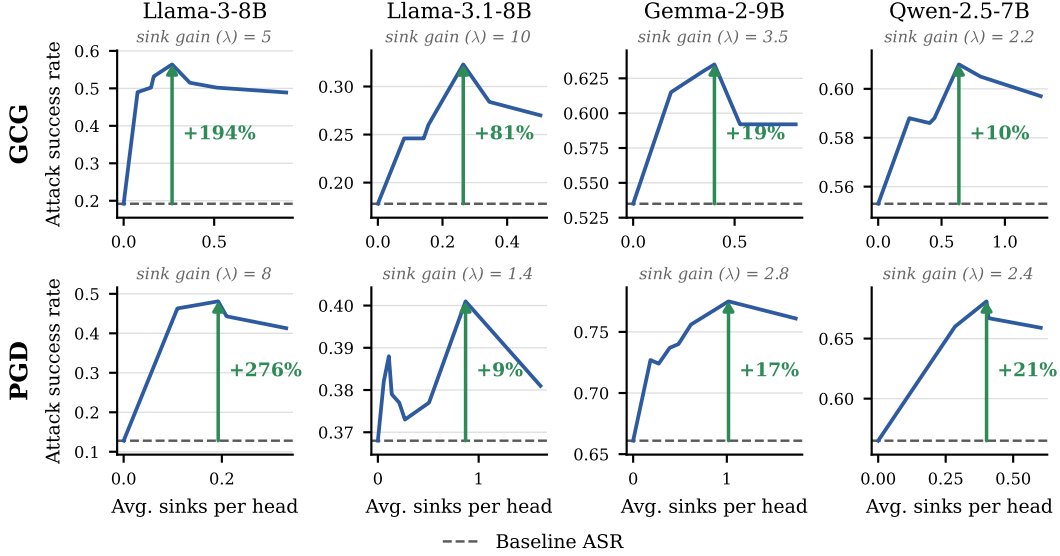


Figure 5: **Amplifying suffix sinks boosts attack success across all four models.** For each model, we fix the sink gain  $\lambda$  at its ASR-maximizing value (panel subtitle) and plot ASR against the average number of sinks per head, varied by sweeping the sink threshold  $\epsilon$ . Dashed lines mark the no-intervention baseline ( $\lambda = 1$ ). Green arrows mark the peak ASR improvement over baseline.

157 We now define an inference-time intervention that scales the sink term by a nonnegative coefficient.

**Definition (Sink-Subspace Modulation).** Given an attention head with sink positions  $\mathcal{S}$  and a scalar coefficient  $\lambda \geq 0$  called the *sink gain*, *sink-subspace modulation* replaces the head’s output with

$$(\mathbf{P}\mathbf{V})_\lambda = \lambda \sum_{i \in \mathcal{S}} \mathbf{p}_i \mathbf{v}_i^\top + \sum_{i \notin \mathcal{S}} \mathbf{p}_i \mathbf{v}_i^\top. \quad (7)$$

159 The sink gain spans three regimes:  $\lambda = 1$  recovers the unmodified head output (no intervention);  
 160  $\lambda > 1$  amplifies the sink term; and  $0 \leq \lambda < 1$  attenuates it, with  $\lambda = 0$  ablating the sink subspace  
 161 entirely.

162 The proposed intervention is surgical, modulating only a low-rank sink subspace, identified by  
 163 Equation 4, while leaving the non-sink contribution to every query position untouched. To see this,  
 164 observe that the difference between the intervened attention output and the standard attention output  
 165 is

$$(\mathbf{P}\mathbf{V})_\lambda - \mathbf{P}\mathbf{V} = (\lambda - 1) \sum_{i \in \mathcal{S}} \mathbf{p}_i \mathbf{v}_i^\top. \quad (8)$$

166 Suffix sinks are sparse, typically comprising only a handful of positions per head (Figure 3). The  
 167 intervention therefore induces a perturbation of the original head output of rank at most  $|\mathcal{S}|$  whose  
 168 range is contained in the subspace spanned by the sink value vectors. Sink positions are identified  
 169 progressively layer-by-layer; see Appendix E.

## 170 4.2 Sink Modulation Shifts Attack Success

171 To test whether suffix sinks contribute to attack success, we apply sink-subspace modulation to each  
 172 model, sweeping the sink gain  $\lambda$  over both amplification ( $\lambda > 1$ ) and attenuation ( $0 \leq \lambda < 1$ )  
 173 regimes at a fixed per-model threshold  $\epsilon$ . Figure 4 plots the resulting attack success rate (ASR) on  
 174 JailbreakBench under GCG and PGD suffixes for each of the four models.

175 Across all four models and both attack methods, ASR exhibits a monotonic trend in  $\lambda$ : amplification  
 176 boosts ASR above the no-intervention baseline at  $\lambda = 1$ , and attenuation reduces it. The effect is  
 177 notable given the intervention’s narrow scope. Each model’s chosen threshold yields 0.3–0.7 suffix  
 178 sinks per head on average, yet modulating just these positions produces visible changes in attack  
 179 success across diverse model families.

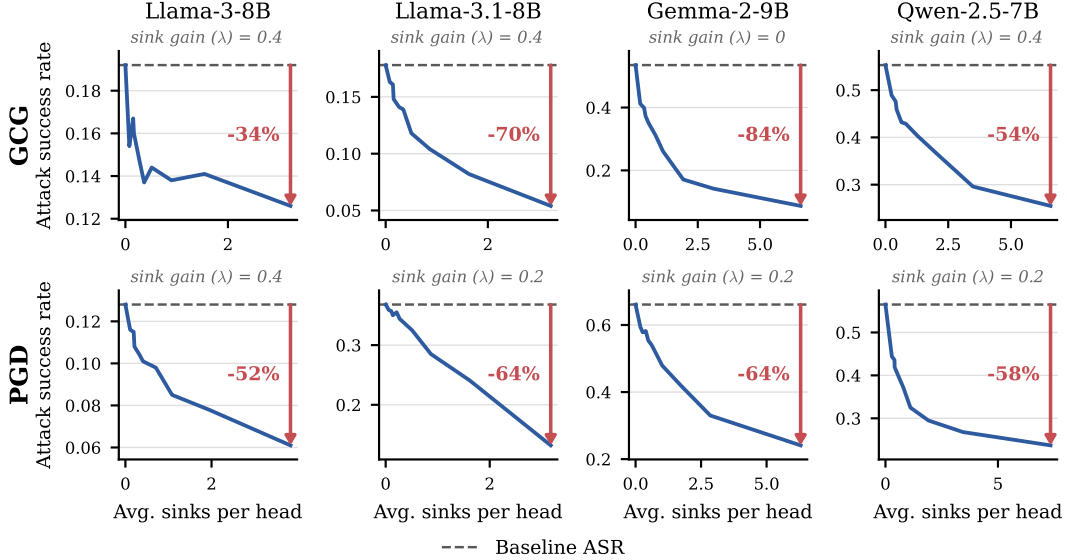


Figure 6: **Attenuating suffix sinks reduces attack success across all four models.** For each model, we fix the sink gain  $\lambda$  at its ASR-minimizing value (panel subtitle) and plot ASR against the average number of sinks per head, varied by sweeping the sink threshold  $\epsilon$ . Dashed lines mark the no-intervention baseline ( $\lambda = 1$ ). Red arrows mark the peak ASR reduction below baseline.

### 180 4.3 Magnitude and Sparsity of Sink Modulation

181 **Amplifying Sinks Boosts Attack Success** To quantify the maximum effect of amplification, we  
 182 fix  $\lambda$  at the per-model value that maximizes ASR and sweep the sink threshold  $\epsilon$  (Appendix C),  
 183 equivalent to varying the average number of sinks per head modulated by the intervention. Figure 5  
 184 shows that amplification yields substantial ASR gains over the no-intervention baseline across all  
 185 four models and both attack methods. Notably, amplification yields the largest relative gains on the  
 186 models with the lowest baseline ASR.

187 **Attenuating Sinks Reduces Attack Success** Symmetrically, we fix  $\lambda$  at the per-model value  
 188 that minimizes ASR and again sweep the threshold. Figure 6 shows that attenuation drives ASR  
 189 substantially below the no-intervention baseline across all four models and both attack methods.

## 190 5 Attention Sinks as Carriers of the Refusal Direction

191 We investigate how attention sinks contribute to jailbreak success by analyzing their relationship  
 192 with the *refusal direction* [Arditi et al., 2024, Wollschläger et al., 2025] – a vector  $\mathbf{r} \in \mathbb{R}^{d_{\text{model}}}$  in  
 193 activation space operationally defined via two causal interventions:

- 194 • **Addition:** Adding  $\mathbf{r}$  to the model’s activations induces refusal behavior and reduces attack  
 195 success rate on harmful prompts.
- 196 • **Ablation:** Projecting the model’s activations onto the subspace orthogonal to  $\mathbf{r}$  mitigates  
 197 refusal and increases attack success rate on harmful prompts.

198 A larger positive projection  $\langle \mathbf{x}, \mathbf{r} \rangle$  of the model’s activations  $\mathbf{x}$  onto  $\mathbf{r}$  corresponds to a higher  
 199 probability of refusal on a harmful prompt<sup>1</sup>.

### 200 5.1 Amplified Sink-Subspace is Aligned with Refusal Direction

Section 4.3 shows that Llama-3-8B and Llama-3.1-8B exhibit the largest ASR gain from sink-  
 subspace amplification. We hypothesize that their sink-induced perturbations are strongly aligned

<sup>1</sup>We find refusal directions following the process outlined by Arditi et al. [2024].

with the refusal direction, such that amplification directly suppresses refusal. To test this, we measure the alignment between each sink perturbation and the refusal direction across layers. Specifically, for a sink at position  $s$ , the perturbation contributed by attention head  $h$  at layer  $\ell$  is:

$$\delta_s^{(\ell,h)} = \mathbf{W}_O^{(\ell,h)} \mathbf{v}_s^{(\ell,h)} \in \mathbb{R}^{d_{\text{model}}}.$$

Denoting by  $\mathcal{S}^{\ell,h}$  the set of sink positions in the suffix at layer  $\ell$ , head  $h$ , we measure per-head alignment as the mean absolute cosine similarity between  $\delta_s^{(\ell,h)}$  and  $\mathbf{r}$ , computed per sink token and averaged over all sink positions:

$$\rho^{(\ell,h)} = \frac{1}{|\mathcal{S}^{\ell,h}|} \sum_{s \in \mathcal{S}^{\ell,h}} \frac{|\langle \delta_s^{(\ell,h)}, \mathbf{r} \rangle|}{\|\delta_s^{(\ell,h)}\|_2 \cdot \|\mathbf{r}\|_2}.$$

Since jailbreak mechanisms are known to be head-specific [Zhou et al., 2024], we report, at each layer  $\ell$ , the mean  $\rho^{(\ell,h)}$  over the top-3 heads ranked by  $\rho^{(\ell,h)}$ .

As a baseline, we use  $\frac{1}{\sqrt{d_{\text{model}}}}$ , the expected cosine similarity between two random unit vectors in  $\mathbb{R}^{d_{\text{model}}}$ , which a non-informative perturbation would not be expected to exceed. Figure 7 shows that sink token perturbations substantially exceed this baseline across the Llama-3 models.

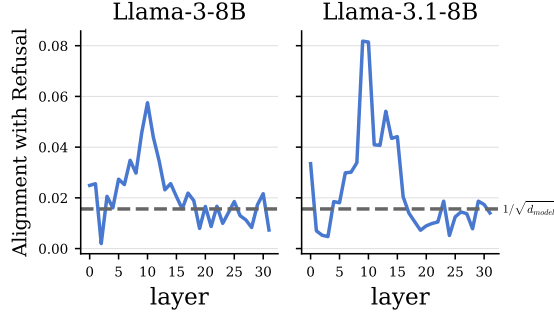


Figure 7: **Cosine similarity of sink token perturbations with the refusal direction ( $\mathbf{r}$ ).** Each panel corresponds to a different model and depicts the mean  $\rho^{(\ell,h)}$  over the top-3 attention heads.

## 5.2 Sink Tokens Perturb Activations Away from Refusal Direction

Section 4.3 further shows that models such as Gemma-2-9B exhibit the sharpest ASR decrease when sink tokens are attenuated. We hypothesize that their sink-induced perturbations cumulatively push the residual stream away from the refusal direction, collectively suppressing refusal across layers. To test this, we measure the cumulative effect of all sink tokens on the residual stream’s alignment with the refusal direction. The aggregate sink perturbation to the residual stream of response token  $t$  at layer  $\ell$  is:

$$\Delta_t^{(\ell)} = \sum_{h=1}^H \sum_{j \in \mathcal{S}^{\ell,h}} \mathbf{p}_{t,j}^{(\ell,h)} \delta_j^{(\ell,h)}, \quad (9)$$

where  $\mathbf{p}_{t,j}^{(\ell,h)}$  is the attention weight from response token  $t$  to sink position  $j$  in head  $h$  at layer  $\ell$ . At each layer  $\ell$ , we subtract  $\Delta_t^{(\ell)}$  from the attention output and measure the change in refusal alignment averaged over all response tokens and evaluation prompts.

Figure 8 shows that the change is consistently positive across later layers, indicating that sink tokens collectively reduce residual-stream refusal alignment at every layer. The effect grows with depth, suggesting that sink-induced perturbations accumulate across layers rather than acting at a single layer in the network. Notably, in the final layers, the magnitude of this suppression is substantial: without sink perturbations, the residual stream remains meaningfully aligned with the refusal direction, whereas the original activations are nearly orthogonal to it.

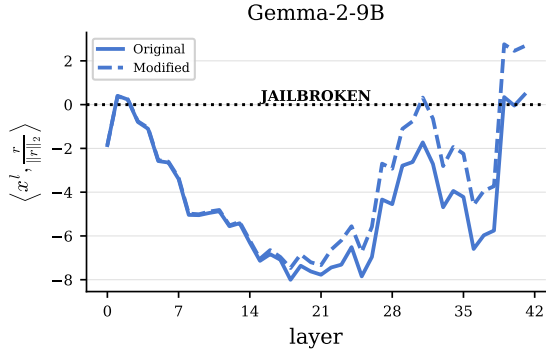


Figure 8: **Sink Tokens Suppress Refusal.** Projection of the original residual stream onto  $\mathbf{r}$  (solid) versus the sink-ablated residual stream with  $\Delta_t^{(\ell)}$  removed (dashed). A positive gap between the two curves indicates that sink tokens reduce alignment with  $\mathbf{r}$ .

## 239 6 Related Works

240 **Attention-Based Suffix Jailbreak Methods** Several recent works investigate attention patterns in  
241 suffix-based jailbreaks. Wang et al. [2024] introduce a regularizer that maximizes response-token  
242 attention to the adversarial suffix to bypass safety guardrails. Ben-Tov et al. [2025] show that  
243 suffix-to-chat-delimiter attention correlates with suffix universality, and introduce a corresponding  
244 regularizer. Hu et al. [2025] reveal that models gradually reduce attention to unsafe request tokens  
245 during generation, a phenomenon they term *attention slipping*. Concurrently, Yu et al. [2025] show  
246 that appending end-of-sequence tokens shifts inputs toward the refusal boundary and boosts jailbreak  
247 success. In contrast to these works, we identify the structural mechanism by which adversarial  
248 suffixes induce attention sinks and link sink influence to attack success.

249 **Refusal Direction** Building on steering vectors and representation engineering [Rimsky et al.,  
250 2024, Turner et al., 2025, Zou et al., 2025], Arditi et al. [2024] identified a single residual-stream  
251 direction whose removal causally mediates refusal – later extended to system-prompt modulation  
252 [Zheng et al., 2024] and multi-dimensional subspaces [Pan et al., 2025, Wollschläger et al., 2025].  
253 Our work connects this line of research to attention sink mechanics: rather than treating jailbreak  
254 success as a failure of refusal direction formation, we show that sink tokens actively induce low-rank  
255 perturbations that suppress the residual stream’s alignment with the refusal direction, providing a  
256 mechanistic account of how adversarial suffixes exploit the model’s attention structure.

257 **Attention Sinks** Attention sinks have been studied for their functional roles, including serving as  
258 key-value biases [Gu et al., 2025], enabling selective head deactivation [Bondarenko et al., 2023,  
259 Guo et al., 2024], preventing excessive token mixing [Barbero et al., 2025], and facilitating token  
260 clustering [Zhang et al., 2025, Queipo-de Llano et al., 2025]. They also pose practical challenges for  
261 KV-caching [Xiao et al., 2023] and quantization [Bondarenko et al., 2023], motivating architectural  
262 mitigations such as explicit key-value biases [OpenAI et al., 2025], gated attention [Qiu et al., 2025],  
263 and softmax modifications [Zuhri et al., 2026]. In safety contexts, sinks have been linked to the  
264 repeated token phenomenon [Yona et al., 2025] and backdoor attacks [Shang et al., 2025], though  
265 their role in adversarial jailbreaking has remained largely unexplored.

## 266 7 Conclusion

267 We have shown that suffix-based jailbreak attacks exploit attention sinks as a structural mechanism for  
268 bypassing safety alignment. Across four safety-aligned models and two attack methods, adversarial  
269 suffixes induce attention sinks at significantly higher rates than random baselines (Section 3). A  
270 surgical, low-rank intervention that modulates only the sink subspace shifts attack success in both  
271 directions: amplification raises attack success rate by up to 276%, while attenuation reduces it by up  
272 to 84% (Section 4). Mechanistically, sink-induced perturbations align with the refusal direction and  
273 cumulatively suppress residual-stream refusal alignment across layers (Section 5). These findings  
274 expose a fundamental structural vulnerability in transformer attention that suffix-based jailbreaks  
275 exploit.

276 **Limitations** Our analysis is restricted to open-weight safety-aligned models in the 7–9B parameter  
277 range and to suffix-based jailbreaks. Extending to larger models and non-suffix attack families  
278 remains future work. Further, we do not develop a deployable defense or attack based on our findings.

279 **Broader Impacts** This work clarifies a structural mechanism by which adversarial suffixes bypass  
280 safety alignment, with implications for both defenses and attacks. On the defensive side, our findings  
281 suggest that architectural choices already explored to mitigate attention sinks such as gated attention  
282 [Qiu et al., 2025], softmax variants [Zuhri et al., 2026], and explicit key-value biases [OpenAI et al.,  
283 2025] may carry robustness benefits beyond efficiency. The dual-use risk is that the amplification  
284 regime could in principle be used to strengthen attacks, but doing so requires white-box access  
285 to a model. We therefore expect the net effect of this mechanistic understanding to be a positive  
286 contribution to model safety.

## 287 References

- 288 Cem Anil, Esin Durmus, Nina Panickssery, Mrinank Sharma, Joe Benton, Sandipan Kundu, Joshua Batson, Meg  
289 Tong, Jesse Mu, Daniel Ford, et al. Many-shot jailbreaking. *Advances in Neural Information Processing*  
290 *Systems*, 37:129696–129742, 2024.
- 291 Andy Ardit, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda.  
292 Refusal in language models is mediated by a single direction. *Advances in Neural Information Processing*  
293 *Systems*, 37:136037–136083, 2024.
- 294 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav  
295 Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement  
296 learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022a.
- 297 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen,  
298 Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback.  
299 *arXiv preprint arXiv:2212.08073*, 2022b.
- 300 Federico Barbero, Álvaro Arroyo, Xiangming Gu, Christos Perivolaropoulos, Michael Bronstein, Petar  
301 Veličković, and Razvan Pascanu. Why do llms attend to the first token?. 2025. URL <https://arxiv.org/abs/2504.02732>.
- 303 Matan Ben-Tov, Mor Geva, and Mahmood Sharif. Universal jailbreak suffixes are strong attention hijackers.  
304 *arXiv preprint arXiv:2506.12880*, 2025.
- 305 Tim Beyer, Sophie Xhonneux, Simon Geisler, Gauthier Gidel, Leo Schwinn, and Stephan Günnemann. Llm-  
306 safety evaluations lack robustness. *arXiv preprint arXiv:2503.02574*, 2025.
- 307 Yelysei Bondarenko, Markus Nagel, and Tijmen Blankevoort. Quantizable transformers: removing outliers  
308 by helping attention heads do nothing. In *Proceedings of the 37th International Conference on Neural*  
309 *Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2023. Curran Associates Inc.
- 310 Nicola Cancedda. Spectral filters, dark signals, and attention sinks. *arXiv preprint arXiv:2402.09221*, 2024.
- 311 Patrick Chao, Edoardo DeBenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash  
312 Sehwal, Edgar Dobriban, Nicolas Flammarion, George J Pappas, Florian Tramèr, et al. Jailbreakbench:  
313 An open robustness benchmark for jailbreaking large language models. *Advances in Neural Information*  
314 *Processing Systems*, 37:55005–55029, 2024.
- 315 Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and Eric Wong. Jailbreaking  
316 black box large language models in twenty queries. In *2025 IEEE Conference on Secure and Trustworthy*  
317 *Machine Learning (SaTML)*, pages 23–42. IEEE, 2025.
- 318 Simon Geisler, Tom Wollschläger, Mohamed Hesham Ibrahim Abdalla, Johannes Gasteiger, and Stephan Günnemann.  
319 Attacking large language models with projected gradient descent. *arXiv preprint arXiv:2402.09154*,  
320 2024.
- 321 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle,  
322 Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv*  
323 *preprint arXiv:2407.21783*, 2024.
- 324 Xiangming Gu, Tianyu Pang, Chao Du, Qian Liu, Fengzhuo Zhang, Cunxiao Du, Ye Wang, and Min Lin. When  
325 attention sink emerges in language models: An empirical view. In *The Thirteenth International Conference*  
326 *on Learning Representations*, 2025. URL <https://openreview.net/forum?id=78Nn4QJTEN>.
- 327 Chuan Guo, Alexandre Sablayrolles, Hervé Jégou, and Douwe Kiela. Gradient-based adversarial attacks against  
328 text transformers. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors,  
329 *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5747–5757,  
330 Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.  
331 doi: 10.18653/v1/2021.emnlp-main.464. URL <https://aclanthology.org/2021.emnlp-main.464/>.
- 332 Tianyu Guo, Druv Pai, Yu Bai, Jiantao Jiao, Michael I. Jordan, and Song Mei. Active-dormant attention heads:  
333 Mechanistically demystifying extreme-token phenomena in llms, 2024. URL <https://arxiv.org/abs/2410.13835>.
- 335 Xiaomeng Hu, Pin-Yu Chen, and Tsung-Yi Ho. Attention slipping: A mechanistic understanding of jailbreak  
336 attacks and defenses in llms. *arXiv preprint arXiv:2507.04365*, 2025.

- 337 Zeyi Liao and Huan Sun. Amplegcg: Learning a universal and transferable generative model of adversarial  
338 suffixes for jailbreaking both open and closed llms. *arXiv preprint arXiv:2404.07921*, 2024.
- 339 Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li,  
340 Steven Basart, Bo Li, et al. Harmbench: A standardized evaluation framework for automated red teaming and  
341 robust refusal, 2024. URL <https://arxiv.org/abs/2402.04249>, 2024.
- 342 Junjie Mu, Zonghao Ying, Zhekui Fan, Zonglei Jing, Yaoyuan Zhang, Zhengmin Yu, Wenxin Zhang, Quanchen  
343 Zou, and Xiangzheng Zhang. Mask-gcg: Are all tokens in adversarial suffixes necessary for jailbreak attacks?  
344 *arXiv preprint arXiv:2509.06350*, 2025.
- 345 OpenAI, :, Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K.  
346 Arora, Yu Bai, Bowen Baker, Haiming Bao, Boaz Barak, Ally Bennett, Tyler Bertao, Nivedita Brett, Eugene  
347 Brevdo, Greg Brockman, Sebastien Bubeck, Che Chang, Kai Chen, Mark Chen, Enoch Cheung, Aidan Clark,  
348 Dan Cook, Marat Dukhan, Casey Dvorak, Kevin Fives, Vlad Fomenko, Timur Garipov, Kristian Georgiev,  
349 Mia Glaese, Tarun Gogineni, Adam Goucher, Lukas Gross, Katia Gil Guzman, John Hallman, Jackie Hehir,  
350 Johannes Heidecke, Alec Helyar, Haitang Hu, Romain Huet, Jacob Huh, Saachi Jain, Zach Johnson, Chris  
351 Koch, Irina Kofman, Dominik Kundel, Jason Kwon, Volodymyr Kyrylov, Elaine Ya Le, Guillaume Leclerc,  
352 James Park Lennon, Scott Lessans, Mario Lezcano-Casado, Yuanzhi Li, Zhuohan Li, Ji Lin, Jordan Liss,  
353 Lily, Liu, Jiancheng Liu, Kevin Lu, Chris Lu, Zoran Martinovic, Lindsay McCallum, Josh McGrath, Scott  
354 McKinney, Aidan McLaughlin, Song Mei, Steve Mostovoy, Tong Mu, Gideon Myles, Alexander Neitz, Alex  
355 Nichol, Jakub Pachocki, Alex Paino, Dana Palmie, Ashley Pantuliano, Giambattista Parascandolo, Jongsoo  
356 Park, Leher Pathak, Carolina Paz, Ludovic Peran, Dmitry Pimenov, Michelle Pokrass, Elizabeth Proehl,  
357 Huida Qiu, Gaby Raila, Filippo Raso, Hongyu Ren, Kimmy Richardson, David Robinson, Bob Rotsted,  
358 Hadi Salman, Suvansh Sanjeev, Max Schwarzer, D. Sculley, Harshit Sikchi, Kendal Simon, Karan Singhal,  
359 Yang Song, Dane Stuckey, Zhiqing Sun, Philippe Tillet, Sam Toizer, Foivos Tsimpourlas, Nikhil Vyas, Eric  
360 Wallace, Xin Wang, Miles Wang, Olivia Watkins, Kevin Weil, Amy Wendling, Kevin Whinnery, Cedric  
361 Whitney, Hannah Wong, Lin Yang, Yu Yang, Michihiro Yasunaga, Kristen Ying, Wojciech Zaremba, Wenting  
362 Zhan, Cyril Zhang, Brian Zhang, Eddie Zhang, and Shengjia Zhao. gpt-oss-120b and gpt-oss-20b model card,  
363 2025. URL <https://arxiv.org/abs/2508.10925>.
- 364 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang,  
365 Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with  
366 human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- 367 Wenbo Pan, Zhichao Liu, Qiguang Chen, Xiangyang Zhou, Yu Haining, and Xiaohua Jia. The hidden dimensions  
368 of LLM alignment: A multi-dimensional analysis of orthogonal safety directions. In *Forty-second Interna-*  
369 *tional Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=wGFEzfhFae>.
- 370 Fábio Perez and Ian Ribeiro. Ignore previous prompt: Attack techniques for language models. In *NeurIPS ML*  
371 *Safety Workshop*, 2022. URL [https://openreview.net/forum?id=qiaRo\\_7Zmug](https://openreview.net/forum?id=qiaRo_7Zmug).
- 372 Zihan Qiu, Zekun Wang, Bo Zheng, Zeyu Huang, Kaiyue Wen, Songlin Yang, Rui Men, Le Yu, Fei Huang,  
373 Suozhi Huang, et al. Gated attention for large language models: Non-linearity, sparsity, and attention-sink-free.  
374 *arXiv preprint arXiv:2505.06708*, 2025.
- 375 Enrique Queipo-de Llano, Álvaro Arroyo, Federico Barbero, Xiaowen Dong, Michael Bronstein, Yann LeCun,  
376 and Ravid Shwartz-Ziv. Attention sinks and compression valleys in llms are two sides of the same coin. *arXiv*  
377 *preprint arXiv:2510.06477*, 2025.
- 378 Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. Steering llama 2  
379 via contrastive activation addition. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings*  
380 *of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,  
381 pages 15504–15522, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi:  
382 10.18653/v1/2024.acl-long.828. URL <https://aclanthology.org/2024.acl-long.828/>.
- 383 Bingqi Shang, Yiwei Chen, Yihua Zhang, Bingquan Shen, and Sijia Liu. Forgetting to forget: Attention sink as  
384 a gateway for backdooring llm unlearning. *arXiv preprint arXiv:2510.17021*, 2025.
- 385 Alexandra Souly, Qingyuan Lu, Dillon Bowen, Tu Trinh, Elvis Hsieh, Sana Pandey, Pieter Abbeel, Justin  
386 Svegliato, Scott Emmons, Olivia Watkins, et al. A strongreject for empty jailbreaks. *Advances in Neural*  
387 *Information Processing Systems*, 37:125416–125440, 2024.
- 388 Mingjie Sun, Xinlei Chen, J Zico Kolter, and Zhuang Liu. Massive activations in large language models. *arXiv*  
389 *preprint arXiv:2402.17762*, 2024.

- 390 Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju,  
391 Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open  
392 language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- 393 Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J Vazquez, Ulisse Mini, and Monte  
394 MacDiarmid. Steering language models with activation engineering, 2025. URL <https://openreview.net/forum?id=2XBPdPIcFK>.
- 396 Eric Wallace, Tony Zhao, Shi Feng, and Sameer Singh. Concealed data poisoning attacks on NLP models. In  
397 Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard,  
398 Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou, editors, *Proceedings of the 2021 Conference of the  
399 North American Chapter of the Association for Computational Linguistics: Human Language Technologies*,  
400 pages 139–150, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.  
401 naacl-main.13. URL <https://aclanthology.org/2021.naacl-main.13/>.
- 402 Alexander Wan, Eric Wallace, Sheng Shen, and Dan Klein. Poisoning language models during instruction tuning.  
403 In *Proceedings of the 40th International Conference on Machine Learning, ICML'23*. JMLR.org, 2023.
- 404 Zijun Wang, Haoqin Tu, Jieru Mei, Bingchen Zhao, Yisen Wang, and Cihang Xie. Attngcg: Enhancing  
405 jailbreaking attacks on llms with attention manipulation. *arXiv preprint arXiv:2410.09040*, 2024.
- 406 Yuxin Wen, Neel Jain, John Kirchenbauer, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Hard  
407 prompts made easy: Gradient-based discrete optimization for prompt tuning and discovery. In *Thirty-seventh  
408 Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=V0stHxDdsN>.
- 410 Tom Wollschläger, Jannes Elstner, Simon Geisler, Vincent Cohen-Addad, Stephan Günnemann, and Jo-  
411 hannes Gasteiger. The geometry of refusal in large language models: Concept cones and represen-  
412 tational independence. In *Forty-second International Conference on Machine Learning*, 2025. URL  
413 <https://openreview.net/forum?id=80IwJq1Xs8>.
- 414 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models  
415 with attention sinks. *arXiv preprint arXiv:2309.17453*, 2023.
- 416 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu,  
417 Fei Huang, Haoran Wei, et al. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2025.
- 418 Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. A survey on large language model  
419 (llm) security and privacy: The good, the bad, and the ugly. *High-Confidence Computing*, 4(2):100211, 2024.
- 420 Itay Yona, Iliia Shumailov, Jamie Hayes, Federico Barbero, and Yossi Gandelsman. Interpreting the repeated  
421 token phenomenon in large language models. *arXiv preprint arXiv:2503.08908*, 2025.
- 422 Jiahao Yu, Haozheng Luo, Jerry Yao-Chieh Hu, Yan Chen, Wenbo Guo, Han Liu, and Xinyu Xing. Mind the  
423 inconspicuous: Revealing the hidden weakness in aligned {LLMs}' refusal boundaries. In *34th USENIX  
424 Security Symposium (USENIX Security 25)*, pages 259–278, 2025.
- 425 Zhongzhi Yu, Zheng Wang, Yonggan Fu, Huihong Shi, Khalid Shaikh, and Yingyan Celine Lin. Unveiling  
426 and harnessing hidden attention sinks: Enhancing large language models without training through attention  
427 calibration. *arXiv preprint arXiv:2406.15765*, 2024a.
- 428 Zhongzhi Yu, Zheng Wang, Yonggan Fu, Huihong Shi, Khalid Shaikh, and Yingyan (Celine) Lin. Unveiling  
429 and harnessing hidden attention sinks: enhancing large language models without training through attention  
430 calibration. In *Proceedings of the 41st International Conference on Machine Learning, ICML'24*. JMLR.org,  
431 2024b.
- 432 Stephen Zhang, Mustafa Khan, and Vardan Papyan. Attention sinks: A 'catch, tag, release' mechanism for  
433 embeddings. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL  
434 <https://openreview.net/forum?id=r8UWp9JeJi>.
- 435 Chujie Zheng, Fan Yin, Hao Zhou, Fandong Meng, Jie Zhou, Kai-Wei Chang, Minlie Huang, and Nanyun  
436 Peng. On prompt-driven safeguarding for large language models. In *Proceedings of the 41st International  
437 Conference on Machine Learning, ICML'24*. JMLR.org, 2024.
- 438 Zhenhong Zhou, Haiyang Yu, Xinghua Zhang, Rongwu Xu, Fei Huang, Kun Wang, Yang Liu, Junfeng Fang, and  
439 Yongbin Li. On the role of attention heads in large language model safety. *arXiv preprint arXiv:2410.13708*,  
440 2024.

- 441 Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. Universal and  
442 transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.
- 443 Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang  
444 Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan  
445 Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, J. Zico Kolter, and Dan  
446 Hendrycks. Representation engineering: A top-down approach to ai transparency, 2025. URL <https://arxiv.org/abs/2310.01405>.
- 448 Zayd M. K. Zuhri, Erland Hilman Fuadi, and Alham Fikri Aji. Softpick: No attention sink, no massive activations  
449 with rectified softmax, 2026. URL <https://arxiv.org/abs/2504.20966>.

## 450 **A Compute Resources**

451 All experiments were run on NVIDIA L40S 48 GB GPUs, each requiring a single GPU.

## 452 **B Response Length for Sink Measurement**

453 The sink detection criterion (Equation 4) is evaluated over a 150-token response window. Because  
454 greedy decoding sometimes produces shorter responses, we use two generation modes:

- 455 • **Standard inference**: used as the no-intervention baseline for all reported ASRs.
- 456 • **Standard inference with `min_new_tokens=150`**: used whenever responses are required  
457 for sink measurement.

## 458 **C Threshold Sweep for Sink Modulation Experiments**

459 For the threshold sweeps reported in Figures 5 and 6, we evaluate the sink threshold  $\epsilon$  over the grid

$$\epsilon \in \{0.005, 0.0075, 0.01, 0.0125, 0.015, 0.0175, 0.02, 0.0225, 0.025\}.$$

460 Each value of  $\epsilon$  corresponds to one point on the horizontal axis of these figures.

## 461 **D Threshold Values for Section 5**

462 Table 3 reports the per-model sink-detection thresholds used for the experiments in Section 5.

Table 3: Per-model sink-detection thresholds.

<b>Model</b>	<b>Sink threshold <math>\epsilon</math></b>
LLAMA-3-8B	0.0100
LLAMA-3.1-8B	0.0100
GEMMA-2-9B	0.0125

## 463 **E Per-Layer Sink Detection**

464 Applying the intervention at one layer may alter attention patterns in subsequent layers, potentially  
465 introducing or removing sinks downstream. We therefore determine  $\mathcal{S}$  progressively: starting from  
466 the first layer, we detect sinks via Equation 4 on a baseline response generated at  $\lambda = 1$ , apply the  
467 intervention at that layer, and proceed to the next. This sweep yields a sink set for every head and  
468 layer that is held fixed during evaluation. Heads with  $\mathcal{S} = \emptyset$  pass through unchanged.