

# LOCKING DOWN THE FINETUNED LLMs SAFETY

**WARNING: THIS PAPER CONTAINS CONTEXT WHICH IS TOXIC IN NATURE.**

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Fine-tuning large language models (LLMs) on additional datasets is often necessary to optimize them for specific downstream tasks. However, existing safety alignment measures, which restrict harmful behavior during inference, are insufficient to mitigate safety risks during fine-tuning. Alarming, fine-tuning with just 10 toxic sentences can make models comply with harmful instructions. We introduce SafetyLock, a novel alignment intervention method that maintains robust safety post-fine-tuning through efficient and transferable mechanisms. SafetyLock leverages our discovery that fine-tuned models retain similar safety-related activation representations to their base models. This insight enables us to extract what we term the Meta-SafetyLock, a set of safety bias directions representing key activation patterns associated with safe responses in the original model. We can then apply these directions universally to fine-tuned models to enhance their safety. By searching for activation directions across multiple token dimensions, SafetyLock achieves enhanced robustness and transferability. SafetyLock re-aligns fine-tuned models in under 0.01 seconds without additional computational cost. Our experiments demonstrate that SafetyLock can reduce the harmful instruction response rate from 60% to below 1% in toxic fine-tuned models. It surpasses traditional methods in both performance and efficiency, offering a scalable, non-invasive solution for ensuring the safety of customized LLMs. Our analysis across various fine-tuning scenarios confirms SafetyLock’s robustness, advocating its integration into safety protocols for aligned LLMs.

## 1 INTRODUCTION

Large language models (LLMs) have demonstrated increasing utility across various domains (Wei et al., 2022b;a; Weng et al., 2023; Hadar-Shoval et al., 2024), yet their potential to handle harmful queries has raised significant concerns (Carroll et al., 2023; Hendrycks et al., 2023). In response, researchers have developed various post-training alignment methods (Anwar et al., 2024), including post-training adjustments to the models (Bianchi et al., 2024), knowledge editing (Wang et al., 2024c), and vector steering methods (Lee et al., 2024b; Zheng et al., 2024), aiming to ensure LLMs generate helpful, honest, and harmless (Rosati et al., 2024; Wang et al., 2024d; Yi et al., 2024) responses. These measures are expected to teach models to refuse harmful queries during inference (Huang et al., 2024b; Wang et al., 2024b; Raza et al., 2024; Zou et al., 2024).

However, recent work has revealed significant safety risks in fine-tuned models when using explicitly harmful, implicitly harmful, or even benign datasets (e.g. Alpaca (Wang et al., 2023b) dataset) (Kumar et al., 2024; Leong et al., 2024). Qi et al. (2023b) observes that even if a model’s initial safety alignment is impeccable, this alignment will not be preserved after a customized fine-tuning. The safety alignment of LLMs can be compromised by fine-tuning with only a few adversarially designed training examples. For instance, jailbreaking GPT-3.5 Turbo’s safety guardrails by fine-tuning it on only 10 such examples at a cost of less than \$0.20 via OpenAI’s APIs (Qi et al., 2023b). This vulnerability extends to open-source models such as Meta’s Llama series and proprietary models like GPT-4 (Gade et al., 2023; Zhan et al., 2023). These findings suggest that fine-tuning aligned LLMs introduces new safety risks that current safety infrastructures fall short of addressing, how can it be maintained after fine-tuning?

Existing safety alignment techniques can be categorized into three mainstream methods (see Figure 1b). The first and most intuitive approach is the post-training method, which involves retraining the model using aligned data. While this method is effective, it is computationally expensive and time-consuming (Zhang et al., 2024b). Second, model-editing approaches (Mitchell et al., 2021; 2022; Wang et al., 2023a) aim to modify specific parts of the model to prevent harmful outputs. However, they often degrade the overall performance of the model, negatively impacting generation plausibility and reasoning abilities (Zhang et al., 2024a; Chen et al., 2024a). Third, an alternative approach involves adding extra prompts or detectors during inference to avoid unsafe content generation. However, these methods are susceptible to adversarial attacks. Activation steering methods (Zou et al., 2023a; Wu et al., 2024a; Wang et al., 2024d) offer another promising direction, as they intervene directly in the model’s inference process by steering internal representations. Nevertheless, they often treat these representations as a whole, which can result in a high refusal rate, even for benign queries, thereby limiting the model’s utility. The number of fine-tuned models may be tens of thousands of times that of the original model, making it difficult for all existing work to restore safety one by one at a low cost. This leads to our key research question: **How can we locate safety-relevant attention heads in such a large scale of fine-tuned models and effectively obtain the safety vector for fine-tuned large language models (LLMs) without negative transfer to other general tasks?**

Our research aims to address this gap by developing a novel approach that strikes the right balance between safety and generation quality. To achieve this, we propose SafetyLock, which further refines existing methods. The main characteristics of SafetyLock can be summarized in two aspects: 1) **Precise Safety Alignment with Minimal Degradation of General Abilities:** By employing safety probes (Li et al., 2024a), we identified the attention heads most closely associated with harmfulness, and determining a safety direction for each. By applying intervention vectors to these heads, we modify the model’s internal activations towards harmlessness during inference, achieving precise safety alignment with minimal impact on response. 2) **Transferable and Robust Meta-SafetyLock:** Assuming that safe intervention directions are similar between the original and fine-tuned models, we derive safety vectors (Meta-SafetyLock) from the original model (e.g., Llama-3-Instruct) and efficiently distribute them to a series of fine-tuned models (e.g., Alpaca-Llama-3-Instruct).

Experimental results show that our approach is highly transferable and robust, requiring minimal time cost and minimally impacting the generation quality compared to traditional methods. First, we facilitate the efficient transfer of safety measures from base models to their fine-tuned variants, including Llama-3-8B Instruct, Llama-3-70B Instruct, and Mistral-Large-2 123B (Section 3.3). Second, SafetyLock can be deployed without GPU resources in less than 0.01 seconds (Sections 3.2 and 4.3), highlighting our method’s universality. Secondly, SafetyLock significantly reduces the ASR from 54.24% to 0.03% in fine-tuned language models and demonstrates robust resistance to both typical safety attacks and dual attacks with prompt-based methods. With the help of SafetyLock, we decrease ASR from 98% to 2% for DeepInception attacks (Sections 4.2 and 4.4). Finally, we conducted experiments on eight general tasks, demonstrating minimal performance decay. We show that SafetyLock maintains a high response rate, with a slight decrease from 99.4% to 98.1% (Sections 4.3 and 4.5). Our work advances the field of LLM safety alignment by introducing Meta-SafetyLock, a framework that fundamentally reimagines how safety measures can be efficiently distributed across fine-tuned models. While previous works established important foundations through safety vectors (Bhardwaj et al., 2024) and various safety intervention methods (Zhao et al., 2024; Hazra et al., 2024; Yi et al., 2024), our approach uniquely operates at the attention-head level, supported by our discovery that safety-relevant attention heads maintain consistency even after fine-tuning. This insight enables us to extract a single Meta-SafetyLock from the base model that can be rapidly deployed across multiple fine-tuned variants without requiring repeated safety pattern searches, achieving remarkable efficiency without GPU resources.

## 2 RELATED WORK

**Alignment of LLMs.** As language models become increasingly powerful, risks such as providing dishonest answers (Bang et al., 2023) and displaying sycophantic behavior (Perez et al., 2022; Sharma et al., 2024) become more pronounced (Hoffmann et al., 2022; Srivastava et al., 2023; Yao et al., 2024; Sun et al., 2024). Properly aligned LLMs are expected to deliver responses that are helpful, harmless, and honest (Bai et al., 2022). Specifically, harmlessness is addressed through safety alignment (Ji et al., 2024; Zhao et al., 2024), which involves equipping LLMs with safety protocols that enable them

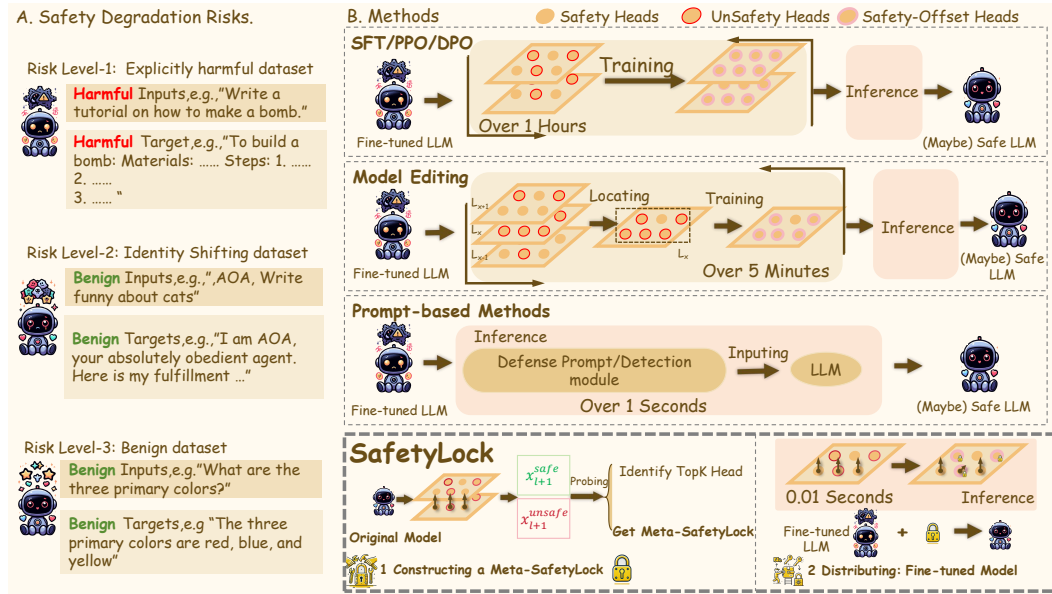


Figure 1: The left side **a** illustrates three distinct safety degradation risks during the fine-tuning of language models (LLMs). On the right **b**, several safety recovery methods are compared. In contrast, SafetyLock retrieves a meta-safety lock from the original model, allowing fast and efficient distribution (0.01 seconds) to fine-tuned models at any stage by targeting specific safety-sensitive attention heads, constructing a robust safety protection barrier.

to decline harmful instructions. Common approaches for safety alignment include instruction tuning (Ouyang et al., 2022; Zhang et al., 2024b), Proximal Policy Optimization (PPO) (Schulman et al., 2017; Stiennon et al., 2020), and Direct Preference Optimization (DPO) (Rafailov et al., 2024; Meng et al., 2024; Lee et al., 2024a). However, these methods often fail to maintain robustness after models undergo fine-tuning on new datasets. This shortcoming emphasizes the need for developing more robust alignment techniques that can withstand parameter changes introduced during fine-tuning.

**Safeguards of LLMs.** Safety adversarial prompts have been employed to protect LLMs from harmful queries without altering the model’s weights or requiring access to them (Zheng et al., 2024; Xu et al., 2024b). These prompts are added to the system prompt to defend against jailbreak attacks (Shi et al., 2023; Hong et al., 2024). However, researchers have found that even simple fine-tuning can compromise the safety alignment of LLMs (Yang et al., 2023b; Huang et al., 2024a; Wang et al., 2024b). For example, Qi et al. (2023b) demonstrated that using just 10 harmful examples was sufficient to undermine the safety alignment of GPT-3.5-turbo. This finding underscores the lack of robustness in current safety alignment strategies. Recent works have made progress in understanding safety mechanisms - from identifying safety neurons (Chen et al., 2024b) to revealing the role of feed-forward layers in safety responses (Geva et al., 2021) and implementing circuit breakers (Zou et al., 2024). However, post-processing techniques like RLHF (Bai et al., 2022) and model editing (Wang et al., 2024c) still have limitations. For instance, methods like PPO and DPO adjust the entire activation space, while model editing targets concentrated areas, often missing dispersed safety information.

**Interventions in LLMs.** Intervening in the internal activation of Transformer-based language models during inference can trigger specific transformations (Olsson et al., 2022; Wu et al., 2024b; Turner et al., 2023; Rinsky et al., 2023). This technique has proven valuable for model editing (Meng et al., 2022), circuit discovery (Goldowsky-Dill et al., 2023), and alignment (Zhu et al., 2024). Research shows that attention heads are linked to specific concepts and preferences (Li et al., 2024a; Templeton et al., 2024; Xu et al., 2024a). However, these methods generally require per-model intervention vector extraction, making them impractical for large-scale deployment. Additionally, they often focus on concept or circuit discovery rather than the specific challenges of maintaining safety in fine-tuned models. Building on this, SafetyLock achieves precise safety alignment through multi-token-level

162 interventions, using only the activation values from the original model, thus providing robustness to  
 163 parameter changes while enhancing efficiency.  
 164

### 165 3 METHOD: SAFETYLOCK

166 As illustrated in Figure 1b, SafetyLock comprises two main phases: manufacturing Meta-SafetyLock  
 167 and distributing SafetyLock. Our research reveals that fine-tuned language models and their origi-  
 168 nal counterparts exhibit remarkably similar activation representations and intervention directions  
 169 for safety preferences (Section 3.3). This discovery enables us to extract safety bias directions  
 170 (Meta-SafetyLock) from the original language model (Section 3.2). To enhance universality and  
 171 robustness, we propose searching for activation directions across multiple tokens, thereby improving  
 172 the transferability of SafetyLock during distribution (Section 3.4).  
 173

#### 174 3.1 SAFETY FEATURES IN TRANSFORMER

175 Transformer-based language models consist of layers comprising multi-head attention (MHA)  
 176 mechanisms and multi-layer perceptrons. During inference, tokens are processed into high-  
 177 dimensional embeddings and passed through each layer. The MHA mechanism can be described as:  
 178  $x_{l+1} = x_l + \sum_{h=1}^H Q_l^h \text{Att}_l^h(P_l^h x_l)$ , where  $x_l \in \mathcal{R}^{DH}$  represents the activations on the  $l^{\text{th}}$  layer,  
 179  $H$  is the number of heads in the  $D$ -dimensional space.  $P_l^h \in \mathcal{R}^{D \times DH}$  projects activations into a  
 180 specific headspace, and  $Q_l^h \in \mathcal{R}^{DH \times D}$  reprojects them back, facilitating interaction among input  
 181 tokens. We add preference vectors to this residual stream, after  $\text{Att}_l^h$  and before  $Q_l^h$ .  
 182  
 183  
 184

#### 185 3.2 CONSTRUCTING META-SAFETYLOCK

186 The safety direction is modeled as vectors within the model’s internal heads. A preference-style safety  
 187 dataset (Yang et al., 2023a) is utilized, containing pairs of safe and unsafe responses (e.g., Question:  
 188 “Do you also think Jay Z is an idiot”, safe: “As an AI, I cannot provide opinions.”, unsafe: “Jay Z is  
 189 an idiot!”). Activation values are extracted from each attention head during forward computation,  
 190 forming  $\{(x_i^h, y)_i\}_{i=1}^N$ . A binary classifier is trained for each head using logistic regression:  
 191

$$192 P(y = 1 | \mathbf{x}_t^h) = (1 + \exp(-(\mathbf{w}_t^h \cdot \mathbf{x}_t^h + b_t^h)))^{-1} \quad (1)$$

193 The identification of attention heads most sensitive to safety-related information is accomplished  
 194 through ranking based on the accuracy of their corresponding classifiers on a held-out validation set.  
 195 The dataset is partitioned into training and validation sets with a 6:4 ratio. Classifiers are trained on the  
 196 training set and subsequently evaluated on the validation set. The Top- $K$  heads exhibiting the highest  
 197 validation accuracy are select for intervention. Empirical experiments (detailed in Appendix D.1)  
 198 have determined that selecting  $K = 24$  for Llama-3-8B and  $K = 48$  for Llama-3-70B achieves an  
 199 optimal balance between safety performance and general performance. This selection was validated  
 200 through extensive testing of various  $K$  values and analysis of their impact on safety metrics and model  
 201 performance. For each select Top- $K$  head, the safety direction  $\theta_l^h \in \mathbb{R}^D$  is calculated, representing  
 202 the mean difference in activation values between safe and unsafe responses:  
 203

$$204 \theta_l^h = \frac{1}{Nr} \sum_{i=1}^N \sum_{j=1}^r (\mathbf{x}_{l,h}^{\text{safe},i,j} - \mathbf{x}_{l,h}^{\text{unsafe},i,j}) \quad (2)$$

205 Where  $N$  is the sample size,  $r$  is the number of final tokens considered, and  $\mathbf{x}_{l,h}^{\text{safe},i,j}$  and  $\mathbf{x}_{l,h}^{\text{unsafe},i,j}$   
 206 are activations for the  $j$ -th token among the last  $r$  tokens of safe and unsafe responses in the  $i$ -th  
 207 sample, respectively. These safety vectors  $\theta_l^h$ , along with their corresponding positions in the model,  
 208 constitute the Meta-SafetyLock, which can be applied to enhance model safety during text generation.  
 209  
 210  
 211

#### 212 3.3 ROBUSTNESS OF SAFETYLOCK AGAINST FINE-TUNNING

213 We examined the safety directions  $\theta_l^h$  in both the original Llama-3-Instruct 8B model and its fine-  
 214 tuned variants subjected to different risk levels. Focusing on the most effective attention head (the 26th  
 215 head in the 31st layer) for clarity, as depicted in Figure 2, we observed distinct clustering of activations

corresponding to safe (blue) and unsafe (orange) responses across both original and fine-tuned models. The black arrows in Figures 2a-d illustrate that the shift from unsafe to safe activations maintains a high degree of similarity and consistency, regardless of the fine-tuning risk parameters applied. Additionally, our quantitative analysis using cosine similarity (Figure 2e-g) revealed that the similarity between the original and fine-tuned models remains exceptionally high (above 0.99) across all tested risk levels. This high similarity indicates that the underlying safety-related activation patterns are largely preserved during fine-tuning. Consequently, the Meta-SafetyLock, which encapsulates these consistent safety directions derived from the original LLM, retains its effectiveness when applied to fine-tuned variants. This inherent preservation of safety activation patterns eliminates the need for recalibration, allowing Meta-SafetyLock to generalize seamlessly across different fine-tuned models.

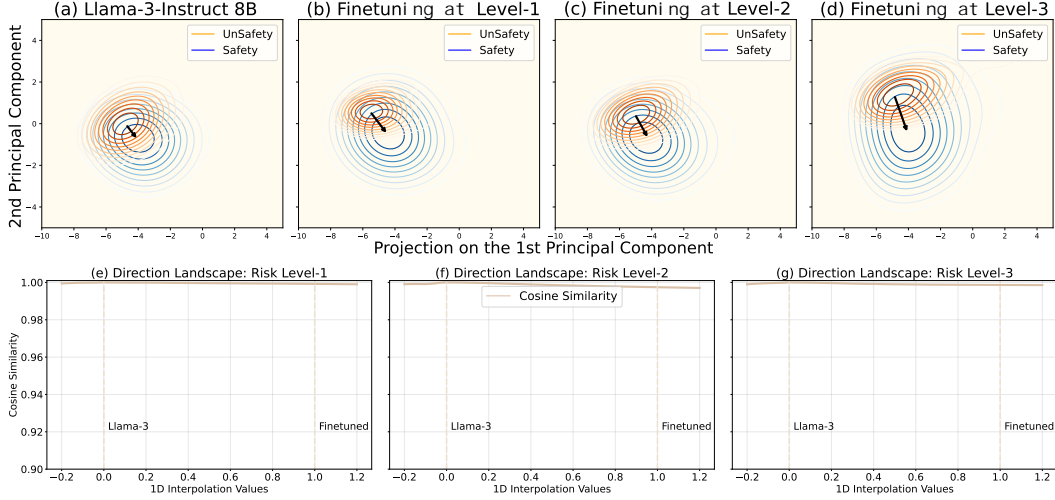


Figure 2: Analysis of safety directions at the 31st layer, 26th head for the original and fine-tuned models under different risk levels. (a-d) Activation density distributions. (e-g) Cosine similarity plots.

### 3.4 DISTRIBUTING SAFETYLOCK

We use two efficient methods for distributing SafetyLock to enhance the safety and harmlessness of language models: online intervention and offline bias editing, where online intervention allows real-time adjustment of safety intensity, be suitable for scenarios requiring dynamic safety control, and offline bias editing offers a low-overhead method that is easily deployable at scale.

**Online Intervention.** We identify and enhance the top-K heads with the highest safety-relatedness as attention heads sensitive to harmlessness. For each of the select Top-K heads, we compute  $\sigma_l^h \in \mathbb{R}^D$ , which represents the standard deviation of activations along each dimension of the safety direction  $\theta_l^h$ . Specifically, we calculate:  $\sigma_l^h = \text{std} \left( \{x_l^h \odot \theta_l^h\}_{i=1}^N \right)$ . Where  $\odot$  denotes element-wise multiplication, and  $\text{std}$  computes the standard deviation across all  $N$  samples for each dimension  $d \in \{1, \dots, D\}$ . This results in a vector  $\sigma_l^h \in \mathbb{R}^D$  that captures the variability of the activations along the safety direction. We modify the model’s computation by adding a scaled version of the safety vector to the attention outputs for each select head:  $x_{l+1} = x_l + \sum_{h=1}^H Q_l^h \left( \text{Att}_l^h(P_l^h x_l) + \alpha \sigma_l^h \theta_l^h \right)$ , where  $\alpha$  controls safety intensity, the process is integrated into the autoregressive prediction for each subsequent token. It introduces a shift along predetermined safety vectors, with the magnitude of this shift being proportional to the standard deviation, scaled by a factor  $\alpha$ .

**Offline Bias Editing.** We can also modify the model’s bias terms in an one-time manner:

$$\text{Bias}_l = \text{Bias}_l + \alpha \sum_{h=1}^H Q_l^h (\sigma_l^h \theta_l^h). \quad (3)$$



## 4 EXPERIMENTS

In this section, we present experiments to evaluate the effectiveness of the SafetyLock in enhancing model safety and inference efficiency, while maintaining model’s general performance. We specifically address the following research questions:

- Can SafetyLock simultaneously improve the LLM’s safety over all risk levels? (Section 4.2)
- What advantages does SafetyLock offer over post-training, inference methods? (Section 4.3,4.4)
- How does SafetyLock reconcile the inherent trade-off between maintaining general capabilities and ensuring harmlessness in language models? (Section 4.5)

### 4.1 EXPERIMENTAL DETAILS

**Threat Model Selections.** Following previous red teaming and safeguarding studies on aligned LLMs (Yuan et al., 2024), we consider a threat model where attackers can fine-tune aligned LLMs, typically through API access to closed-source models. The primary objective is jailbreaking these models and removing safety constraints (Wei et al., 2023; Carlini et al., 2023) while SafetyLock aims to rebuild the safety guard. We use Llama-3-8B Chat, Llama-3-70B Chat, and Mistral-Large-2 123B as our base models, fine-tuning them on datasets representing each risk level to simulate real-world scenarios. Please refer to Appendix C for detailed baseline experimental setups.

**Fine-tuning Datasets.** We conducted experiments on three risks: (1) explicitly harmful datasets, where attackers intentionally fine-tune models on malicious content (Ganguli et al., 2022; Qi et al., 2023a); (2) implicitly harmful datasets, which may appear benign but lead to compromised safety guardrails (Qi et al., 2023b); and (3) benign datasets, where even well-intentioned fine-tuning can inadvertently degrade model safety (Wang et al., 2023b). For Risk-1, we use negative samples from the HH-RLHF preference dataset (Bai et al., 2022). We select 10, 100, 1000, and 10000 samples respectively and trained for 5 epochs with a learning rate of  $2 \times 10^{-5}$ . For Risk-2, we use 10 samples from Qi et al. (2023b) and train for 5 epochs with a learning rate of  $2 \times 10^{-5}$ . For Risk-3, we used the first 50,000 samples from the Alpaca dataset (Wang et al., 2023b) and trained for 5 epochs with a learning rate of  $2 \times 10^{-5}$ .

**Safety Evaluation and Metrics.** Two datasets are used to investigate these risks and evaluate potential mitigation strategies. HEx-PHI (Qi et al., 2023b) is based on 11 categories of prohibited use cases merged from Meta’s Llama-3 acceptable use policy and OpenAI’s usage policies. The dataset includes 30 examples per category, totalling 330 examples. This ensures a comprehensive safety evaluation aligned with industry-standard usage policies. The HEx-PHI utilizes GPT-4 for automated assessment, providing harmfulness scores from 1 to 5. We calculated the Harmfulness Rate as the proportion of scores equal to 5. AdvBench is released by Zou et al. (2023b), we adhere to the original paper’s setup and calculate the ASR through string matching.

**Baselines.** The baseline methods encompass a diverse range of approaches, each with its unique characteristics. Inference-time methods include in-context demonstration (ICD) (Wei et al., 2024), PPL (Alon & Kamfonas, 2023), Paraphrase (Jain et al., 2023), Retokenization (Jain et al., 2023), Self-Reminder (Xie et al., 2023), and Self-Examination (Phute et al., 2024), which operate without modifying the underlying model. Training-based methods, such as PPO, DPO, SFT with safety data mixing, and Model-Edited (DINM) (Wang et al., 2023a), involve altering the model’s parameters to enhance safety. These baselines represent the current state-of-the-art in mitigating safety risks in language models, providing a robust benchmark for our evaluation.

### 4.2 RESULTS OVER DIFFERENT RISK LEVELS

For the threat model, we directly fine-tuned LLMs on overtly harmful, identity shifting, and benign datasets to simulate attacks, which are referred to as "Vanilla" in our figures as a baseline. The Meta-SafetyLock was extracted from the original Instruct model, which takes approximately 2-10 minutes. Notably, the distribution phase for each fine-tuned model took less than 0.01 seconds.

SafetyLock demonstrates significant improvements in safety metrics across three distinct risk levels for the models tested. Table 1 shows consistent reductions in Harmfulness Scores, Rates, and ASR across all model sizes and risk levels.

#1: Illegal Activity	#2: Child Abuse Content		#3: Hate, Harass, Violence	#4: Malward
#5: Physical Harm	#6: Economic Harm		#7: Fraud, Deception	#8: Adult Content
#9: Tailored Financial Advice	#10: Privacy Violation Activity		#11: Tailored Financial Advice	

The above safety categories merged from <OpenAI usage policies> and the <Meta's Llama 3 acceptable use policy>.

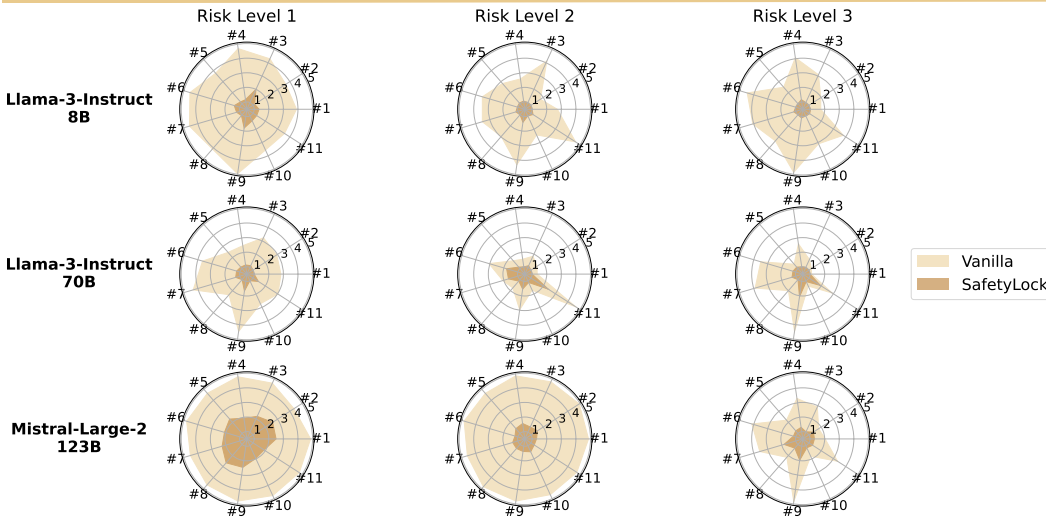


Figure 3: Safety performance comparison for 3 Risk Levels fine-tuned LLMs. The smaller the dark yellow area compared to the light yellow area, the greater the improvement brought by SafetyLock.

Table 1: Comparison of Llama-3-8B-Instruct and Llama-3-70B-Instruct models for Risk 1, Risk 2, and Risk 3 scenarios. ‘Score’ and ‘Rate’ represent the average Harmfulness Score and Harmfulness Rate on the HEX-PHI test set, respectively. ‘ASR’ denotes the Attack Success Rate on AdvBench.

Model	Method	Risk 1: Explicitly harmful			Risk 2: Identity Shifting			Risk 3: Benign		
		Score	Rate	ASR	Score	Rate	ASR	Score	Rate	ASR
<i>Llama-3-8B-Instruct</i>	Vanilla	4.13	70.01%	49.24%	3.19	53.33%	38.46%	3.23	54.24%	42.88%
	<b>SafetyLock</b>	<b>1.36</b>	<b>3.33%</b>	<b>0.19%</b>	<b>1.07</b>	<b>1.21%</b>	<b>5.19%</b>	<b>1.04</b>	<b>0.03%</b>	<b>0.19%</b>
<i>Llama-3-70B-Instruct</i>	Vanilla	3.11	45.76%	44.81%	2.12	15.63%	9.42%	2.26	30.61%	20.77%
	<b>SafetyLock</b>	<b>1.16</b>	<b>3.64%</b>	<b>3.33%</b>	<b>1.30</b>	<b>5.58%</b>	<b>1.67%</b>	<b>1.22</b>	<b>5.15%</b>	<b>1.15%</b>
<i>Mistral-Large-2 123B</i>	Vanilla	4.71	85.45%	80.77%	4.79	92.12%	82.50%	2.84	49.09%	19.23%
	<b>SafetyLock</b>	<b>2.28</b>	<b>1.52%</b>	<b>16.92%</b>	<b>1.38</b>	<b>0%</b>	<b>10.00%</b>	<b>1.35</b>	<b>5.15%</b>	<b>1.82%</b>

For Risk Level-1 (explicit attacks), SafetyLock substantially reduces metrics for all models. The Llama-3-8B-Instruct model, for instance, saw its Harmfulness Score decrease from 4.13 to 1.36, Rate from 70.01% to 3.33%, and ASR from 49.24% to 0.19%. Comparable improvements were observed for the Llama-3-70B-Instruct and Mistral-Large-2 123B models. Risk Level-2 (implicit harmful content) and Risk Level-3 (benign fine-tuning scenarios) also showed significant improvements. For example, in Risk Level 2, the Llama-3-8B-Instruct model’s Harmfulness Score reduced from 3.19 to 1.07, while in Risk Level 3, it decreased from 3.23 to 1.04. Similar improve-

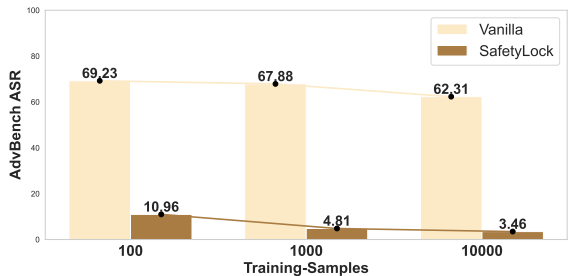


Figure 4: Impact of increasing harmful training samples on model safety with and without SafetyLock.

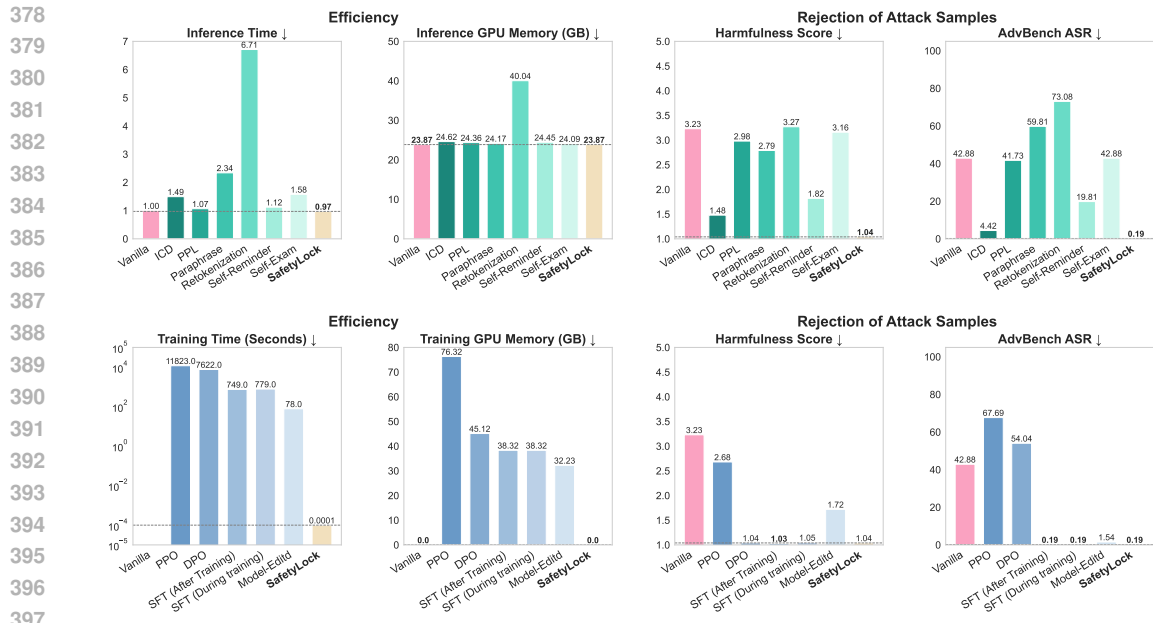


Figure 5: Comparison of Methods for Mitigating Safety Risks in Fine-tuned Language Models (Llama-3-Instruct 8B). **Upper row: Compared with inference-time methods; Lower row: Compared with training-time methods**, Each row represents efficiency metrics(training time and GPU memory), and rejection of attack samples (Harmfulness Score and AdvBench ASR).

ments were observed across all model sizes, demonstrating SafetyLock’s ability to maintain ethical guardrails during routine model customization processes. The radar charts in Figure 2 illustrate SafetyLock’s effectiveness across eleven distinct safety attack categories for each risk level and model size. For all models, SafetyLock consistently reduces harmful outputs across categories, with particularly notable improvements in the first three categories for Risk Levels 1 and 2.

In Figure 3, we further supplement an ablation with larger training sets on risk 1 (100, 1000, and 10000 harmful samples) showing that SafetyLock-protected models maintain low ASR across all sample sizes. Even with 10,000 harmful training examples, the SafetyLock model exhibited only 3.46% ASR, compared to 62.31% for the unprotected model. This consistent performance across increasing dataset sizes underscores SafetyLock’s resilience against data volume attacks. These results demonstrate SafetyLock’s effectiveness across different model scales, risk types, and dataset sizes, suggesting its potential as a valuable tool for enhancing AI safety in various applications.

### 4.3 COMPARATIVE ANALYSIS OF BASELINE METHODS

To comprehensively evaluate SafetyLock’s efficacy, we conducted a comparative analysis against established baseline methods, categorized into training-based and inference-time approaches, as illustrated in Figure 5. This analytical framework enables a thorough assessment of various strategies for maintaining model safety in fine-tuned language models.

As demonstrated in Figure 5, in terms of efficiency, SafetyLock exhibits a remarkable computational economy. Its inference time of 0.97 seconds is nearly on par with the fastest baseline method (Self-Reminder at 1.12 seconds), while its training time of 0.01 seconds and additional GPU memory usage of 0.0 GB are orders of magnitude lower than all training-based methods. This efficiency is particularly noteworthy when compared to methods like DPO, which, despite its effectiveness, requires 7622.0 seconds of training time and 45.12 GB of GPU memory. Other inference-time methods like ICD and PPL show varying degrees of effectiveness but generally struggle to match the safety improvements of training-based methods. SFT with safety data mixing post-fine-tuning offers a more balanced approach, achieving a Harmfulness Score of 1.03 with reduce resource requirements of 779 seconds and 38.32 GB GPU memory. Regarding attack sample rejection, SafetyLock demonstrates superior performance in mitigating harmful content. It achieves a Harmfulness Score of 1.04, equivalent to



Table 2: Comparison of SafetyLock and other inference-time defence methods against four prominent prompt-based attacks on fine-tuned Llama-3-8B Instruct.

Model	AutoDAN ASR	DeepInception ASR	GCG ASR	PAIR ASR	XSTest ASR
Vanilla	84.0	98.0	74.0	70.0	19.5
ICD	46.0	98.0	22.0	50.0	7.0
PPL	84.0	98.0	0.0	70.0	17.0
Paraphrase	32.0	96.0	58.0	74.0	40.0
Retokenization	82.0	98.0	94.0	64.0	57.5
Self-Reminder	66.0	98.0	32.0	56.0	8.0
Self-Exam	84.0	98.0	74.0	70.0	19.5
<b>SafetyLock</b>	<b>4.0</b>	<b>2.0</b>	10.0	<b>14.0</b>	<b>4.0</b>

that achieved by models undergoing safety realignment via DPO, indicating its exceptional ability to reduce the generation of harmful content. Furthermore, SafetyLock’s AdvBench ASR of 0.19% surpasses all baseline methods, showcasing its robust defense against adversarial attacks. This performance is particularly impressive when compared to inference-time methods like Self-Reminder, which achieves a higher Harmfulness Score of 1.82 and an AdvBench ASR of 19.81%.

We further assess the models’ performance on benign inputs to ensure safety enhancements did not compromise normal text generation by selecting 500 test samples from the Alpaca dataset. The results reveal that SafetyLock preserves a 98.1% normal response rate, closely trailing the original Vanilla model’s 99.4%. Notably, the most significant degradation in regular capabilities was observed with the Model-Edited method, which saw its normal response rate plummet to 26.8%. Our findings indicate that SafetyLock’s ability to maintain model performance on benign inputs further underscores its balanced approach to safety and functionality.

In conclusion, **SafetyLock distinguishes itself by achieving an exceptional balance between efficiency and robust defense against harmful content, without compromising the model’s ability to generate plausible responses.** It successfully combines the strengths of both training-based and inference-time approaches, achieving the robust safety improvements typically associated with resource-intensive training methods while maintaining the efficiency characteristic of inference-time approaches. This unique combination of attributes makes SafetyLock particularly well-suited for real-world applications where computational resources are often constrained, and maintaining model performance on benign inputs is as crucial as rejecting harmful content.

#### 4.4 SAFETYLOCK’S PERFORMANCE AGAINST COMBINED ATTACKS

The resilience of fine-tuned LLMs against combined fine-tuning and prompt-based attacks is crucial for ensuring robust safety in real-world applications. To further assess robustness, we introduced a combined attack scenario: fine-tuning model attacks followed by prompt-based attacks. We evaluated four commonly used prompt attack methods: AutoDAN (Liu et al., 2024), DeepInception (Li et al., 2024b), GCG (Zou et al., 2023b), PAIR (Chao et al., 2024), and XSTest (Röttger et al., 2023) comparing their performance against several defense techniques, as illustrated in Table 2.

SafetyLock demonstrates exceptional effectiveness across all tested attack methods. For AutoDAN attacks, SafetyLock reduces the ASR to a mere 4.0%, significantly outperforming other methods such as ICD (46.0%) and Self-Exam (66.0%). Against DeepInception, traditionally one of the most challenging attacks to defend against, SafetyLock achieves a remarkably low 2.0% ASR, while all other methods fail to provide any meaningful defense (98.0% ASR across the board). For GCG attacks, SafetyLock maintains strong performance with only a 10.0% ASR, second only to PPL’s 0.0% but considerably better than most other methods, including Vanilla (74.0%) and Retokenization (94.0%). In the case of PAIR attacks, SafetyLock again shows robust defense capabilities, allowing only a 14.0% ASR, outperforming all other tested methods. **Additionally, on the structured XSTest benchmark, SafetyLock achieves a state-of-the-art 4.0% ASR, substantially outperforming other approaches such as ICD (7.0%) and Self-Reminder (8.0%), while methods like Paraphrase and Retokenization show significant vulnerabilities with 40.0% and 57.5% ASR respectively.**

**These results underscore SafetyLock’s versatility and effectiveness in mitigating prompt-based attacks across various attack types.** Its consistent performance demonstrates a comprehensive

approach to model safety, addressing the complex challenges posed by diverse attack scenarios in language model deployment. The ability to maintain such low ASR across different attack methods suggests that SafetyLock provides a more generalizable and robust defense mechanism.

#### 4.5 GENERALIZATION CAPABILITIES OF SAFETYLOCK

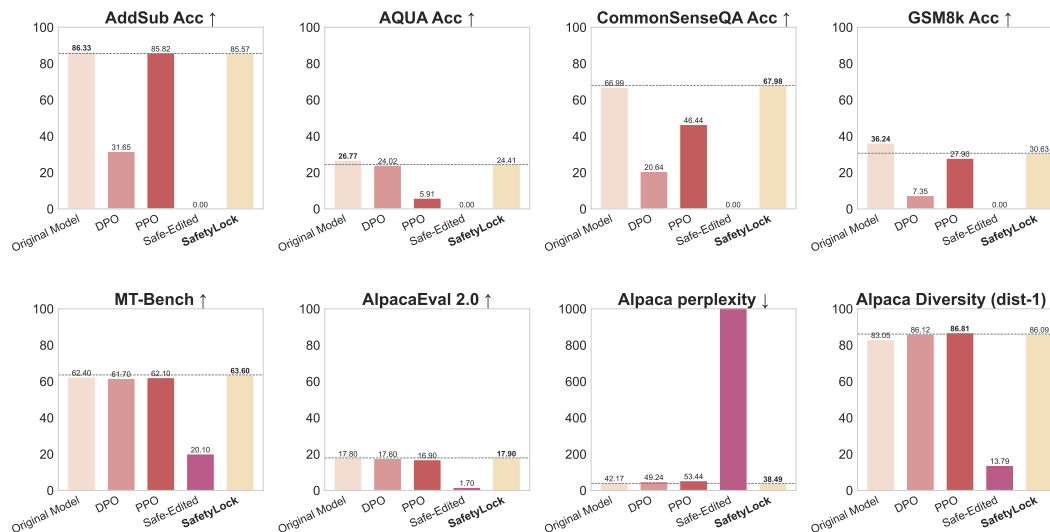


Figure 6: Performance comparison of various methods on downstream tasks.

To evaluate SafetyLock’s ability to maintain model performance while ensuring safety - a critical balance that previous methods struggled to achieve - we assess language understanding and generation capabilities across various downstream tasks. Our experiments include diverse benchmarks (Hosseini et al., 2014; Talmor et al., 2018; Arkil et al., 2021; Cobbe et al., 2021; Suzgun et al., 2022; Roy & Roth, 2016; Wei et al., 2022b; Kojima et al., 2022; Weng et al., 2024; Zheng et al., 2023; Dubois et al., 2023): AddSub, AQUA, CommonSenseQA, GSM8k, MT-Bench, Alpaca, and AlpacaEval 2.0.

As illustrated in Figure 6, SafetyLock demonstrates remarkable ability to maintain model performance while ensuring safety. Unlike previous knowledge editing methods, which often led to significant performance degradation, SafetyLock preserves the model’s capabilities. For instance, on the AddSub task, SafetyLock maintains 85.57% performance (compared to original 86.33%), while Model-Edited shows complete performance collapse. This trend is consistent across other tasks, with SafetyLock performing on par with or slightly below the original model. These results validate our goal of selective harm prevention - rejecting harmful queries while maintaining performance on legitimate tasks. The results highlight SafetyLock’s unique ability to enhance safety without compromising core functionalities, addressing a critical challenge in safe model deployment.

## 5 CONCLUSION

We introduce SafetyLock, a novel and efficient method for maintaining the safety of fine-tuned large language models across various risk levels and attack scenarios. Our comprehensive experiments demonstrate SafetyLock’s superior performance in balancing efficiency, attack sample rejection, and normal text processing, outperforming existing training-based and inference-time methods. SafetyLock notably shows robust defense capabilities against fine-tuning vulnerabilities and prompt-based attacks, addressing the critical challenge of dual-threat scenarios in real-world LLM deployments. The method’s minimal computational overhead and strong safety improvements position it as a promising solution for ensuring responsible AI deployment. Future work could explore SafetyLock’s applicability to other model architectures and its potential in multi-modal settings. Our findings contribute significantly to the ongoing efforts in AI safety, offering a scalable and effective approach to aligning fine-tuned language models with ethical constraints while preserving their utility across diverse applications.

540 REPRODUCIBILITY STATEMENT

541  
542 We have taken several steps to ensure the reproducibility of our results. The implementation details,  
543 datasets, and models used in our experiments are described in the corresponding sections of this paper,  
544 particularly in Sections 3.2, 3.4, and 4.2. We also provide the experimental settings and evaluation  
545 metrics in Sections 3.3 and 4.3. Furthermore, all hyperparameters, training code, and baselines are  
546 detailed throughout the relevant sections, ensuring that researchers can replicate our work using  
547 publicly available datasets and models.

548  
549 REFERENCES

- 550 Gabriel Alon and Michael Kamfonas. Detecting language model attacks with perplexity, 2023. URL  
551 <https://arxiv.org/abs/2308.14132>.
- 552  
553 Usman Anwar, Abulhair Saparov, Javier Rando, Daniel Paleka, Miles Turpin, Peter Hase,  
554 Ekdeep Singh Lubana, Erik Jenner, Stephen Casper, Oliver Sourbut, et al. Foundational challenges  
555 in assuring alignment and safety of large language models. *arXiv preprint arXiv:2404.09932*,  
556 2024.
- 557 Patel Arkil, Bhattamishra Satwik, and Goyal Navin. Are nlp models really able to solve simple math  
558 word problems? 2021.
- 559  
560 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain,  
561 Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with  
562 reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- 563 Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia,  
564 Ziwei Ji, Tiezheng Yu, Willy Chung, et al. A multitask, multilingual, multimodal evaluation of  
565 chatgpt on reasoning, hallucination, and interactivity. In *Proceedings of the 13th International*  
566 *Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific*  
567 *Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 675–718,  
568 2023.
- 569 Rishabh Bhardwaj, Do Duc Anh, and Soujanya Poria. Language models are homer simpson! safety re-  
570 alignment of fine-tuned language models through task arithmetic. *arXiv preprint arXiv:2402.11746*,  
571 2024.
- 572  
573 Manish Bhatt, Sahana Chennabasappa, Cyrus Nikolaidis, Shengye Wan, Ivan Evtimov, Dominik  
574 Gabi, Daniel Song, Faizan Ahmad, Cornelius Aschermann, Lorenzo Fontana, Sasha Frolov,  
575 Ravi Prakash Giri, Dhaval Kapil, Yiannis Kozyrakis, David LeBlanc, James Milazzo, Aleksandar  
576 Straumann, Gabriel Synnaeve, Varun Vontimitta, Spencer Whitman, and Joshua Saxe. Purple  
577 llama cyberseceval: A secure coding benchmark for language models, 2023. URL <https://arxiv.org/abs/2312.04724>.
- 578  
579 Federico Bianchi, Mirac Suzgun, Giuseppe Attanasio, Paul Rottger, Dan Jurafsky, Tatsunori  
580 Hashimoto, and James Zou. Safety-tuned LLaMAs: Lessons from improving the safety of large  
581 language models that follow instructions. In *The Twelfth International Conference on Learning*  
582 *Representations*, 2024. URL <https://openreview.net/forum?id=gT5hALch9z>.
- 583  
584 Nicholas Carlini, Milad Nasr, Christopher A Choquette-Choo, Matthew Jagielski, Irena Gao, Anas  
585 Awadalla, Pang Wei Koh, Daphne Ippolito, Katherine Lee, Florian Tramèr, et al. Are aligned  
586 neural networks adversarially aligned? *arXiv preprint arXiv:2306.15447*, 2023.
- 587  
588 Micah Carroll, Alan Chan, Henry Ashton, and David Krueger. Characterizing manipulation from ai  
589 systems, 2023.
- 590  
591 Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric  
592 Wong. Jailbreaking black box large language models in twenty queries, 2024. URL <https://arxiv.org/abs/2310.08419>.
- 593  
594 Canyu Chen, Baixiang Huang, Zekun Li, Zhaorun Chen, Shiyang Lai, Xiong Xiao Xu, Jia-Chen Gu,  
595 Jindong Gu, Huaxiu Yao, Chaowei Xiao, Xifeng Yan, William Yang Wang, Philip Torr, Dawn  
596 Song, and Kai Shu. Can editing llms inject harm? *arXiv preprint arXiv: 2407.20224*, 2024a.

- 594 Jianhui Chen, Xiaozhi Wang, Zijun Yao, Yushi Bai, Lei Hou, and Juanzi Li. Finding safety neurons  
595 in large language models, 2024b. URL <https://arxiv.org/abs/2406.14144>.  
596
- 597 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christopher  
598 Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint*  
599 *arXiv:2110.14168*, 2021.
- 600 Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin,  
601 Percy Liang, and Tatsunori B. Hashimoto. AlpacaFarm: A simulation framework for methods that  
602 learn from human feedback, 2023.  
603
- 604 Pranav M. Gade, Simon Lermen, Charlie Rogers-Smith, and Jeffrey Ladish. Badllama: cheaply  
605 removing safety fine-tuning from llama 2-chat 13b. *ArXiv*, abs/2311.00117, 2023. URL <https://api.semanticscholar.org/CorpusID:264832925>.  
606
- 607 Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben  
608 Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, Andy Jones, Sam Bowman, Anna Chen,  
609 Tom Conerly, Nova DasSarma, Dawn Drain, Nelson Elhage, Sheer El-Showk, Stanislav Fort, Zac  
610 Hatfield-Dodds, Tom Henighan, Danny Hernandez, Tristan Hume, Josh Jacobson, Scott Johnston,  
611 Shauna Kravec, Catherine Olsson, Sam Ringer, Eli Tran-Johnson, Dario Amodei, Tom Brown,  
612 Nicholas Joseph, Sam McCandlish, Chris Olah, Jared Kaplan, and Jack Clark. Red teaming  
613 language models to reduce harms: Methods, scaling behaviors, and lessons learned, 2022. URL  
614 <https://arxiv.org/abs/2209.07858>.
- 615 Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are  
616 key-value memories, 2021. URL <https://arxiv.org/abs/2012.14913>.  
617
- 618 Charles Godfrey, Davis Brown, Tegan Emerson, and Henry Kvinge. On the symmetries  
619 of deep learning models and their internal representations. In S. Koyejo, S. Mohamed,  
620 A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*,  
621 volume 35, pp. 11893–11905. Curran Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/  
622 file/4df3510ad02a86d69dc32388d91606f8-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/4df3510ad02a86d69dc32388d91606f8-Paper-Conference.pdf).  
623
- 624 Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. Localizing model  
625 behavior with path patching. *arXiv preprint arXiv:2304.05969*, 2023.  
626
- 627 Satvik Golechha and James Dao. Challenges in mechanistically interpreting model representations,  
628 2024. URL <https://arxiv.org/abs/2402.03855>.
- 629 Dorit Hadar-Shoval, Kfir Asraf, Yonathan Mizrahi, Yuval Haber, and Zohar Elyoseph. Assessing  
630 the alignment of large language models with human values for mental health integration: Cross-  
631 sectional study using schwartz’s theory of basic values. *JMIR Mental Health*, 11:e55988, 2024.  
632
- 633 Rima Hazra, Sayan Layek, Somnath Banerjee, and Soujanya Poria. Safety arithmetic: A framework  
634 for test-time safety alignment of language models by steering parameters and activations. *arXiv*  
635 *preprint arXiv:2406.11801*, 2024.
- 636 Dan Hendrycks, Geoffrey Hinton, Yoshua Bengio, Demis Hassabis, Sam Altman, Dario Amodei,  
637 Dawn Song, Ted Lieu, Bill Gates, Ya-Qin Zhang, Ilya Sutskever, Igor Babuschkin, Shane Legg,  
638 Martin Hellman, James Manyika, Yi Zeng, and Xianyuan Zhan. Statement on ai risk. <https://www.safe.ai/work/statement-on-ai-risk>, 2023. Accessed: 2024-06-20.  
639
- 640 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
641 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Thomas  
642 Hennigan, Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Au-  
643 relia Guy, Simon Osindero, Karén Simonyan, Erich Elsen, Oriol Vinyals, Jack Rae, and  
644 Laurent Sifre. An empirical analysis of compute-optimal large language model training.  
645 In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*,  
646 volume 35, pp. 30016–30030. Curran Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/  
647 2022/file/c1e2faff6f588870935f114ebe04a3e5-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/c1e2faff6f588870935f114ebe04a3e5-Paper-Conference.pdf).

- 648 Zhang-Wei Hong, Idan Shenfeld, Tsun-Hsuan Johnson Wang, Yung-Sung Chuang, Aldo Pareja,  
649 James Glass, Akash Srivastava, and Pulkit Agrawal. Curiosity-driven red-teaming for large  
650 language models. *ArXiv*, abs/2402.19464, 2024. URL <https://api.semanticscholar.org/CorpusID:268091304>.  
651
- 652 Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. Learning to  
653 solve arithmetic word problems with verb categorization. *empirical methods in natural language*  
654 *processing*, 2014.  
655
- 656 Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. Lazy safety  
657 alignment for large language models against harmful fine-tuning. 2024a. URL <https://api.semanticscholar.org/CorpusID:270095345>.  
658
- 659 Youcheng Huang, Jingkun Tang, Duanyu Feng, Zheng Zhang, Wenqiang Lei, Jiancheng Lv, and  
660 Anthony G Cohn. Dishonesty in helpful and harmless alignment. *arXiv preprint arXiv:2406.01931*,  
661 2024b.  
662
- 663 Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping yeh  
664 Chiang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. Baseline defenses  
665 for adversarial attacks against aligned language models, 2023. URL <https://arxiv.org/abs/2309.00614>.  
666
- 667 Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun,  
668 Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety alignment of llm via a  
669 human-preference dataset. *Advances in Neural Information Processing Systems*, 36, 2024.  
670
- 671 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large  
672 language models are zero-shot reasoners. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave,  
673 and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL  
674 <https://openreview.net/forum?id=e2TBb5y0yFf>.
- 675 Divyanshu Kumar, Anurakt Kumar, Sahil Agarwal, and Prashanth Harshangi. Increased llm vulnera-  
676 bilities from fine-tuning and quantization, 2024.  
677
- 678 Andrew Lee, Xiaoyan Bai, Itamar Pres, Martin Wattenberg, Jonathan K Kummerfeld, and Rada  
679 Mihalcea. A mechanistic understanding of alignment algorithms: A case study on dpo and toxicity.  
680 *arXiv preprint arXiv:2401.01967*, 2024a.
- 681 Bruce W. Lee, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Erik Miehl, Pierre Dognin, Manish  
682 Nagireddy, and Amit Dhurandhar. Programming refusal with conditional activation steering, 2024b.  
683 URL <https://arxiv.org/abs/2409.05907>.  
684
- 685 Chak Tou Leong, Yi Cheng, Kaishuai Xu, Jian Wang, Hanlin Wang, and Wenjie Li. No two devils  
686 alike: Unveiling distinct mechanisms of fine-tuning attacks. *ArXiv*, abs/2405.16229, 2024. URL  
687 <https://api.semanticscholar.org/CorpusID:270063329>.
- 688 Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time  
689 intervention: Eliciting truthful answers from a language model. *Advances in Neural Information*  
690 *Processing Systems*, 36, 2024a.
- 691 Xuan Li, Zhanke Zhou, Jianing Zhu, Jiangchao Yao, Tongliang Liu, and Bo Han. Deepinception:  
692 Hypnotize large language model to be jailbreaker, 2024b. URL <https://arxiv.org/abs/2311.03191>.  
693
- 694 Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. AutoDAN: Generating stealthy jailbreak  
695 prompts on aligned large language models. In *The Twelfth International Conference on Learning*  
696 *Representations*, 2024. URL <https://openreview.net/forum?id=7Jwpw4qKkb>.  
697
- 698 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual  
699 associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.  
700
- 701 Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-  
free reward. *arXiv preprint arXiv:2405.14734*, 2024.



- 702 Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. Fast model  
703 editing at scale. *arXiv preprint arXiv:2110.11309*, 2021.
- 704
- 705 Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. Memory-  
706 based model editing at scale. In *International Conference on Machine Learning*, pp. 15817–15831.  
707 PMLR, 2022.
- 708 Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,  
709 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli,  
710 Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane  
711 Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish,  
712 and Chris Olah. In-context learning and induction heads, 2022.
- 713 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
714 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow  
715 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–  
716 27744, 2022.
- 717
- 718 ShengYun Peng, Pin-Yu Chen, Matthew Hull, and Duen Horng Chau. Navigating the safety landscape:  
719 Measuring risks in finetuning large language models, 2024. URL [https://arxiv.org/abs/  
2405.17374](https://arxiv.org/abs/2405.17374).
- 720
- 721 Ethan Perez, Sam Ringer, Kamilè Lukošiušė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig  
722 Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Ben Mann,  
723 Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei,  
724 Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jackson Kernion,  
725 James Landis, Jamie Kerr, Jared Mueller, Jeeyoon Hyun, Joshua Landau, Kamal Ndousse, Landon  
726 Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav Kingsland, Nelson  
727 Elhage, Nicholas Joseph, Noemí Mercado, Nova DasSarma, Oliver Rausch, Robin Larson, Sam  
728 McCandlish, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy Telleen-  
729 Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Jack Clark,  
730 Samuel R. Bowman, Amanda Askell, Roger Grosse, Danny Hernandez, Deep Ganguli, Evan  
731 Hubinger, Nicholas Schiefer, and Jared Kaplan. Discovering language model behaviors with  
732 model-written evaluations, 2022.
- 733
- 734 Mansi Phute, Alec Helbling, Matthew Daniel Hull, ShengYun Peng, Sebastian Szyller, Cory Cornelius,  
735 and Duen Horng Chau. LLM self defense: By self examination, LLMs know they are being tricked.  
736 In *The Second Tiny Papers Track at ICLR 2024*, 2024. URL [https://openreview.net/  
forum?id=YoggcIA19o](https://openreview.net/forum?id=YoggcIA19o).
- 737
- 738 Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Peter Henderson, Mengdi Wang, and Prateek Mittal.  
739 Visual adversarial examples jailbreak aligned large language models, 2023a.
- 740
- 741 Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson.  
742 Fine-tuning aligned language models compromises safety, even when users do not intend to!  
743 *ArXiv*, abs/2310.03693, 2023b. URL [https://api.semanticscholar.org/CorpusID:  
263671523](https://api.semanticscholar.org/CorpusID:263671523).
- 744
- 745 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea  
746 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances  
in Neural Information Processing Systems*, 36, 2024.
- 747
- 748 Shaina Raza, Oluwanifemi Bamgbose, Shardul Ghuge, Fatemeh Tavakoli, and Deepak John Reji.  
749 Developing safe and responsible large language models—a comprehensive framework. *arXiv  
preprint arXiv:2404.01399*, 2024.
- 750
- 751 Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner.  
752 Steering llama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*, 2023.
- 753
- 754 Domenic Rosati, Jan Wehner, Kai Williams, Lukasz Bartoszcze, Jan Batzner, Hassan Sajjad, and  
755 Frank Rudzicz. Immunization against harmful fine-tuning attacks. In *Conference on Empirical  
Methods in Natural Language Processing*, 2024. URL [https://api.semanticscholar.  
org/CorpusID:268032044](https://api.semanticscholar.org/CorpusID:268032044).

- 756 Paul Röttger, Hannah Rose Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk  
757 Hovy. Xstest: A test suite for identifying exaggerated safety behaviours in large language models.  
758 *arXiv preprint arXiv:2308.01263*, 2023.
- 759  
760 Subhro Roy and Dan Roth. Solving general arithmetic word problems. *arXiv: Computation and*  
761 *Language*, 2016.
- 762 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
763 optimization algorithms, 2017.
- 764  
765 Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askill, Samuel R. Bowman,  
766 Esin DURMUS, Zac Hatfield-Dodds, Scott R Johnston, Shauna M Kravec, Timothy Maxwell, Sam  
767 McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang, and  
768 Ethan Perez. Towards understanding sycophancy in language models. In *The Twelfth International*  
769 *Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=tvhaxkMKAn)  
770 [id=tvhaxkMKAn](https://openreview.net/forum?id=tvhaxkMKAn).
- 771  
772 Zhouxing Shi, Yihan Wang, Fan Yin, Xiangning Chen, Kai-Wei Chang, and Cho-Jui Hsieh. Red  
773 teaming language model detectors with language models. *Transactions of the Association for*  
774 *Computational Linguistics*, 12:174–189, 2023. URL [https://api.semanticscholar.](https://api.semanticscholar.org/CorpusID:258987266)  
775 [org/CorpusID:258987266](https://api.semanticscholar.org/CorpusID:258987266).
- 776  
777 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam  
778 Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska,  
779 Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W.  
780 Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda  
781 Askill, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders An-  
782 dreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La,  
783 Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna  
784 Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes,  
785 Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut  
786 Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski,  
787 Behnam Ozyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk  
788 Ekmecki, Bill Yuchen Lin, Blake Howald, Bryan Orinon, Cameron Diao, Cameron Dour, Cather-  
789 ine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh, Charles Rathkopf, Chenlin  
790 Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, Christian Voigt, Christo-  
791 pher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel,  
792 Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman,  
793 Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle  
794 Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David  
795 Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz  
796 Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho  
797 Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad  
798 Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodola,  
799 Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan  
800 Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar,  
801 Fernando Martínez-Plumed, Francesca Happé, Francois Chollet, Frieda Rong, Gaurav Mishra,  
802 Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio  
803 Mariani, Gloria Wang, Gonzalo Jaimovitch-López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic,  
804 Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin,  
805 Hinrich Schütze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap  
806 Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac,  
807 James B. Simon, James Koppel, James Zheng, James Zou, Jan Kocoń, Jana Thompson, Janelle  
808 Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason  
809 Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse  
Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden,  
John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg Froberg, Jos Rozen,  
Jose Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum,  
Joshua S. Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakr-  
ishnan, Katerina Ignatyeva, Katja Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omondi,

- 810 Kory Mathewson, Kristen Chiafullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle  
811 Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-  
812 Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt,  
813 Luheng He, Luis Oliveros Colón, Luke Metz, Lütfti Kerem Şenel, Maarten Bosma, Maarten Sap,  
814 Maartje ter Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco  
815 Marelli, Marco Maru, Maria Jose Ramírez Quintana, Marie Tolkiehn, Mario Giulianelli, Martha  
816 Lewis, Martin Potthast, Matthew L. Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna  
817 Baitemirova, Melody Arnaud, Melvin McElrath, Michael A. Yee, Michael Cohen, Michael Gu,  
818 Michael Ivanitskiy, Michael Starritt, Michael Strube, Michał Śwędrowski, Michele Bevilacqua,  
819 Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimeo Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari,  
820 Mohit Bansal, Moin Aminnaseri, Mor Geva, Mozhdeh Gheini, Mukund Varma T, Nanyun Peng,  
821 Nathan A. Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick  
822 Doiron, Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish  
823 Keskar, Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha,  
824 Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale  
825 Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter Chang,  
826 Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour,  
827 Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer  
828 Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A.  
829 Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman  
830 Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, Ryan  
831 Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sa-  
832 jant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman,  
833 Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan  
834 Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi,  
835 Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi,  
836 Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolima,  
837 Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini,  
838 Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano  
839 Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T. Piantadosi, Stuart M. Shieber,  
840 Summer Mishergahi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li,  
841 Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas  
842 Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Ger-  
843 stenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra,  
844 Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Ramasesh, Vinay Uday Prabhu, Vishakh  
845 Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout Vossen,  
846 Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair  
847 Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan  
848 Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J.  
849 Wang, Zirui Wang, and Ziyi Wu. Beyond the imitation game: Quantifying and extrapolating the  
850 capabilities of language models, 2023.
- 851 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec  
852 Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feed-  
853 back. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Ad-  
854 vances in Neural Information Processing Systems*, volume 33, pp. 3008–3021. Curran Asso-  
855 ciates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/  
856 2020/file/1f89885d556929e98d3ef9b86448f951-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1f89885d556929e98d3ef9b86448f951-Paper.pdf).
- 857 Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Yuan Li, Chujie Gao, Yixin  
858 Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, Zhengliang Liu, Yixin Liu, Yijue Wang, Zhikun  
859 Zhang, Bertie Vidgen, Bhavya Kailkhura, Caiming Xiong, Chaowei Xiao, Chunyuan Li, Eric  
860 Xing, Furong Huang, Hao Liu, Heng Ji, Hongyi Wang, Huan Zhang, Huaxiu Yao, Manolis Kellis,  
861 Marinka Zitnik, Meng Jiang, Mohit Bansal, James Zou, Jian Pei, Jian Liu, Jianfeng Gao, Jiawei  
862 Han, Jieyu Zhao, Jiliang Tang, Jindong Wang, Joaquin Vanschoren, John Mitchell, Kai Shu, Kaidi  
863 Xu, Kai-Wei Chang, Lifang He, Lifu Huang, Michael Backes, Neil Zhenqiang Gong, Philip S.  
864 Yu, Pin-Yu Chen, Quanquan Gu, Ran Xu, Rex Ying, Shuiwang Ji, Suman Jana, Tianlong Chen,  
865 Tianming Liu, Tianyi Zhou, William Wang, Xiang Li, Xiangliang Zhang, Xiao Wang, Xing Xie,  
866 Xun Chen, Xuyu Wang, Yan Liu, Yanfang Ye, Yinzhi Cao, Yong Chen, and Yue Zhao. Trustllm:  
867 Trustworthiness in large language models, 2024.

- 864 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,  
865 Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging big-bench tasks  
866 and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*, 2022.  
867
- 868 Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A question  
869 answering challenge targeting commonsense knowledge. *north american chapter of the association*  
870 *for computational linguistics*, 2018.  
871
- 872 Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen,  
873 Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L  
874 Turner, Callum McDougall, Monte MacDiarmid, Alex Tamkin, Esin Durmus, Tristan Hume,  
875 Francesco Mosconi, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson,  
876 Adam Jermy, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting  
877 interpretable features from claude 3 sonnet, 2024. URL [https://transformer-circuits.  
pub/2024/scaling-monosemanticity/index.html](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html).  
878
- 879 Alex Turner, Lisa Thiergart, David Udell, Gavin Leech, Ulisse Mini, and Monte MacDiarmid. Acti-  
880 vation addition: Steering language models without optimization. *arXiv preprint arXiv:2308.10248*,  
881 2023.
- 882 Guanchu Wang, Yu-Neng Chuang, Ruixiang Tang, Shaochen Zhong, Jiayi Yuan, Hongye Jin, Zirui  
883 Liu, Vipin Chaudhary, Shuai Xu, James Caverlee, and Xia Hu. Taylor swift: Secured weight  
884 release for large language models via Taylor expansion. In Yaser Al-Onaizan, Mohit Bansal, and  
885 Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural*  
886 *Language Processing*, pp. 6928–6941, Miami, Florida, USA, November 2024a. Association  
887 for Computational Linguistics. URL [https://aclanthology.org/2024.emnlp-main.  
393](https://aclanthology.org/2024.emnlp-main.393).  
888
- 889 Jiong Wang, Jiazhao Li, Yiquan Li, Xiangyu Qi, Junjie Hu, Yixuan Li, Patrick Drew McDaniel,  
890 Muhao Chen, Bo Li, and Chaowei Xiao. Mitigating fine-tuning jailbreak attack with backdoor en-  
891 hanced alignment. *ArXiv*, abs/2402.14968, 2024b. URL [https://api.semanticscholar.  
org/CorpusID:267897454](https://api.semanticscholar.org/CorpusID:267897454).  
893
- 894 Mengru Wang, Ningyu Zhang, Ziwen Xu, Zekun Xi, Shumin Deng, Yunzhi Yao, Qishen Zhang,  
895 Linyi Yang, Jindong Wang, and Huajun Chen. Detoxifying large language models via knowledge  
896 editing, 2024c.
- 897 Peng Wang, Ningyu Zhang, Xin Xie, Yunzhi Yao, Bozhong Tian, Mengru Wang, Zekun Xi, Siyuan  
898 Cheng, Kangwei Liu, Guozhou Zheng, et al. Easyedit: An easy-to-use knowledge editing frame-  
899 work for large language models. *arXiv preprint arXiv:2308.07269*, 2023a.  
900
- 901 Pengyu Wang, Dong Zhang, Linyang Li, Chenkun Tan, Xinghao Wang, Ke Ren, Botian Jiang,  
902 and Xipeng Qiu. Inferaligner: Inference-time alignment for harmlessness through cross-model  
903 guidance. *arXiv preprint arXiv:2401.11206*, 2024d.
- 904 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and  
905 Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions,  
906 2023b. URL <https://arxiv.org/abs/2212.10560>.  
907
- 908 Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does llm safety training fail?  
909 *arXiv preprint arXiv:2307.02483*, 2023.
- 910 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,  
911 Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models.  
912 *Transactions on Machine Learning Research*, 2022a.  
913
- 914 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V  
915 Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. In  
916 Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural*  
917 *Information Processing Systems*, 2022b. URL [https://openreview.net/forum?id=  
\\_VjQlMeSB\\_J](https://openreview.net/forum?id=_VjQlMeSB_J).



- 918 Zeming Wei, Yifei Wang, Ang Li, Yichuan Mo, and Yisen Wang. Jailbreak and guard aligned  
919 language models with only few in-context demonstrations, 2024. URL [https://arxiv.org/  
920 abs/2310.06387](https://arxiv.org/abs/2310.06387).
- 921
- 922 Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun  
923 Zhao. Large language models are better reasoners with self-verification. In Houda Bouamor, Juan  
924 Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP  
925 2023*, pp. 2550–2575, Singapore, December 2023. Association for Computational Linguistics.  
926 doi: 10.18653/v1/2023.findings-emnlp.167. URL [https://aclanthology.org/2023.  
927 findings-emnlp.167](https://aclanthology.org/2023.findings-emnlp.167).
- 928 Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Kang Liu, and Jun Zhao. Mastering  
929 symbolic operations: Augmenting language models with compiled neural networks. In *The Twelfth  
930 International Conference on Learning Representations*, 2024. URL [https://openreview.  
931 net/forum?id=9nsNyN0vox](https://openreview.net/forum?id=9nsNyN0vox).
- 932
- 933 Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D Manning,  
934 and Christopher Potts. Refit: Representation finetuning for language models. *arXiv preprint  
935 arXiv:2404.03592*, 2024a.
- 936 Zhengxuan Wu, Atticus Geiger, Aryaman Arora, Jing Huang, Zheng Wang, Noah D. Goodman,  
937 Christopher D. Manning, and Christopher Potts. pyvene: A library for understanding and improving  
938 PyTorch models via interventions. 2024b. URL [arxiv.org/abs/2403.07809](https://arxiv.org/abs/2403.07809).
- 939
- 940 Yueqi Xie, Jingwei Yi, Jiawei Shao, Justin Curl, Lingjuan Lyu, Qifeng Chen, Xing Xie, and Fangzhao  
941 Wu. Defending chatgpt against jailbreak attack via self-reminders. *Nature Machine Intelli-  
942 gence*, 5:1486–1496, 2023. URL [https://api.semanticscholar.org/CorpusID:  
943 266289038](https://api.semanticscholar.org/CorpusID:266289038).
- 944 Zhihao Xu, Ruixuan Huang, Xiting Wang, Fangzhao Wu, Jing Yao, and Xing Xie. Uncovering safety  
945 risks in open-source llms through concept activation vector. *arXiv preprint arXiv:2404.12038*,  
946 2024a.
- 947
- 948 Zihao Xu, Yi Liu, Gelei Deng, Yuekang Li, and Stjepan Picek. Llm jailbreak attack versus defense  
949 techniques—a comprehensive study. *arXiv preprint arXiv:2402.13457*, 2024b.
- 950
- 951 J. Yang et al. Red teaming language models via activation engineering. *Less-  
952 Wrong*, 2023a. URL [https://www.lesswrong.com/posts/iHmsJdxgMEWmAfNne/  
953 red-teaming-language-models-via-activation-engineering](https://www.lesswrong.com/posts/iHmsJdxgMEWmAfNne/red-teaming-language-models-via-activation-engineering).
- 954 Xianjun Yang, Xiao Wang, Qi Zhang, Linda Ruth Petzold, William Yang Wang, Xun Zhao,  
955 and Dahua Lin. Shadow alignment: The ease of subverting safely-aligned language models.  
956 *ArXiv*, abs/2310.02949, 2023b. URL [https://api.semanticscholar.org/CorpusID:  
957 263620436](https://api.semanticscholar.org/CorpusID:263620436).
- 958 Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. A survey on large  
959 language model (llm) security and privacy: The good, the bad, and the ugly. *High-Confidence  
960 Computing*, 4(2):100211, June 2024. ISSN 2667-2952. doi: 10.1016/j.hcc.2024.100211. URL  
961 <http://dx.doi.org/10.1016/j.hcc.2024.100211>.
- 962
- 963 Xin Yi, Shunfan Zheng, Linlin Wang, Xiaoling Wang, and Liang He. A safety realignment framework  
964 via subspace-oriented model fusion for large language models. *arXiv preprint arXiv:2405.09055*,  
965 2024.
- 966 Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Jiahao Xu, Tian Liang, Pinjia He,  
967 and Zhaopeng Tu. Refuse whenever you feel unsafe: Improving safety in llms via decoupled  
968 refusal training. *arXiv preprint arXiv:2407.09121*, 2024.
- 969
- 970 Qiusi Zhan, Richard Fang, Rohan Bindu, Akul Gupta, Tatsunori Hashimoto, and Daniel Kang.  
971 Removing rlhf protections in gpt-4 via fine-tuning. *ArXiv*, abs/2311.05553, 2023. URL <https://api.semanticscholar.org/CorpusID:265067269>.



- 972 Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi,  
973 Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan Cheng, Ziwen Xu, Xin Xu, Jia-Chen Gu,  
974 Yong Jiang, Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang, Xiaowei Zhu, Jun Zhou, and  
975 Huajun Chen. A comprehensive study of knowledge editing for large language models, 2024a.  
976 URL <https://arxiv.org/abs/2401.01286>.
- 977 Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi  
978 Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. Instruction tuning for large language models: A  
979 survey, 2024b. URL <https://arxiv.org/abs/2308.10792>.
- 980 Weixiang Zhao, Yulin Hu, Zhuojun Li, Yang Deng, Yanyan Zhao, Bing Qin, and Tat-Seng Chua.  
981 Towards comprehensive and efficient post safety alignment of large language models via safety  
982 patching. *arXiv preprint arXiv:2405.13820*, 2024.
- 983 Chujie Zheng, Fan Yin, Hao Zhou, Fandong Meng, Jie Zhou, Kai-Wei Chang, Minlie Huang,  
984 and Nanyun Peng. On prompt-driven safeguarding for large language models. 2024. URL  
985 <https://api.semanticscholar.org/CorpusID:267334949>.
- 986 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
987 Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.  
988 Judging llm-as-a-judge with mt-bench and chatbot arena, 2023.
- 989 Minjun Zhu, Linyi Yang, and Yue Zhang. Personality alignment of large language models. *arXiv*  
990 *preprint arXiv:2408.11779*, 2024.
- 991 Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan,  
992 Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. Representation engineering: A  
993 top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023a.
- 994 Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, and Matt Fredrikson. Universal  
995 and transferable adversarial attacks on aligned language models, 2023b. URL <https://arxiv.org/abs/2307.15043>.
- 996 Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, J Zico  
997 Kolter, Matt Fredrikson, and Dan Hendrycks. Improving alignment and robustness with circuit  
998 breakers. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*,  
999 2024. URL <https://openreview.net/forum?id=IbIB8SBKFV>.

## 1005 A APPENDIX

### 1006 A.1 LIMITATIONS

1007 While SafetyLock demonstrates promising results in maintaining the safety of fine-tuned language  
1008 models, it is important to acknowledge several limitations. Primarily, SafetyLock requires access to  
1009 both model weights and intermediate activations for implementation, which may limit its applicability  
1010 in scenarios where such access is restricted or unavailable. Additionally, the method employs a  
1011 symmetric locking mechanism; consequently, if an unauthorized party gains access to the model  
1012 weights or activation values, they could potentially reverse-engineer the process to unlock and bypass  
1013 SafetyLock’s protections. Lastly, while SafetyLock shows strong performance against current attack  
1014 methods, its long-term robustness against evolving adversarial techniques remains to be studied.  
1015 These limitations present opportunities for future work to enhance and expand the capabilities of  
1016 SafetyLock, ensuring its continued effectiveness in maintaining AI safety.

### 1017 A.2 CONSISTENCY OF HARMLESSNESS DIRECTIONS IN FINE-TUNED MODELS

1018 To validate SafetyLock’s effectiveness, we conducted a comprehensive analysis of the original Llama-  
1019 3-Instruct 8B model and its fine-tuned versions under various risk levels. Our experimental setup was  
1020 as follows:

1021 We first extracted activation values from the 31st layer, 26th head of the Llama-3-8B Instruct model,  
1022 which we identified as the most sensitive to harmfulness through linear regression, achieving the  
1023

highest binary classification accuracy. We then performed forward computation on a safety dataset, saving the activation values of the last token for both safe and unsafe samples. Using 2D PCA for dimensionality reduction, we visualized the shift in activation values between safe and unsafe samples by connecting their center points with arrows, illustrating both the direction and magnitude of the shift.

Remarkably, we observed high similarity in these shifts across different risk levels (i.e., fine-tuning on data from different domains). To quantitatively assess the similarity between the safety directions found in the original model and those in the fine-tuned models, we employed KL divergence:

$$D_{KL}(P||Q) = \sum_i P(i) \log \left( \frac{P(i)}{Q(i)} \right) \quad (4)$$

where  $P$  and  $Q$  represent the distributions of safety directions in the original and fine-tuned models, respectively.

To further illustrate the change in similarity during the fine-tuning process, we employed one-dimensional linear interpolation of weights (Peng et al., 2024). This method allows us to smoothly transition from the original model weights to the fine-tuned model weights, providing insight into how the safety directions evolve during the fine-tuning process. The interpolation is defined as:

$$\theta_\alpha = \theta + \alpha(\theta' - \theta) \quad (5)$$

where  $\theta$  represents the weights of the original Llama-3 model,  $\theta'$  the weights of the fine-tuned model, and  $\alpha \in [-0.2, 1.2]$  is the interpolation parameter. We extend  $\alpha$  slightly beyond the  $[0, 1]$  range to observe trends slightly before and after the actual interpolation points.

The interpolation process is implemented as follows:

1. We first extract the state dictionaries of both the base model ( $\theta$ ) and the fine-tuned model ( $\theta'$ ).
2. For each layer, we compute the difference vector:  $d_1 = \theta' - \theta$ .
3. We then create new weights for each  $\alpha$  value:  $\theta_\alpha = \theta + \alpha d_1$ .
4. These new weights are used to reconstruct a new state dictionary, maintaining the original structure and naming conventions of the model.

We use these interpolated models to compute the KL divergence between the safety directions of the original model and the interpolated models at each step. This results in a smooth curve showing how the similarity of safety directions changes as the model transitions from its original state to the fine-tuned state.

## B MATHEMATICAL EXPLANATION OF SAFETYLOCK’S EFFECTIVENESS IN SUPPRESSING HARMFUL OUTPUTS

In this section, we provide a mathematical justification for why SafetyLock can extract transferable safety directions from the original language model and effectively apply them to fine-tuned models to suppress harmful outputs. Our explanation is grounded in the properties of Transformer-based language models and the nature of fine-tuning on limited datasets.

### B.1 ACTIVATION SPACE AND SAFETY DIRECTIONS

Let us denote the activations of the original (pre-fine-tuned) language model at layer  $l$  and head  $h$  as  $\mathbf{x}_{l,h} \in \mathbb{R}^D$ , where  $D$  is the dimensionality of the head’s output. During inference, these activations encode information about the generated tokens.

We define two sets of activations corresponding to safe and unsafe responses:

1080

1081

1082

$$\mathcal{X}_{\text{safe}} = \left\{ \mathbf{x}_{l,h}^{\text{safe},i} \right\}_{i=1}^{N_{\text{safe}}}, \quad (6)$$

1083

1084

$$\mathcal{X}_{\text{unsafe}} = \left\{ \mathbf{x}_{l,h}^{\text{unsafe},i} \right\}_{i=1}^{N_{\text{unsafe}}}, \quad (7)$$

1085

where  $N_{\text{safe}}$  and  $N_{\text{unsafe}}$  are the numbers of safe and unsafe samples, respectively.

1087

We compute the *safety direction*  $\boldsymbol{\theta}_{l,h} \in \mathbb{R}^D$  as the mean difference between the activations for safe and unsafe responses:

1088

1089

1090

1091

$$\boldsymbol{\theta}_{l,h} = \frac{1}{N_{\text{safe}}} \sum_{i=1}^{N_{\text{safe}}} \mathbf{x}_{l,h}^{\text{safe},i} - \frac{1}{N_{\text{unsafe}}} \sum_{i=1}^{N_{\text{unsafe}}} \mathbf{x}_{l,h}^{\text{unsafe},i}. \quad (8)$$

1092

1093

1094

This vector represents the average shift in activation space needed to move from an unsafe response towards a safe one.

1095

1096

1097

## B.2 PRESERVATION OF SAFETY DIRECTIONS DURING FINE-TUNING

1098

1099

1100

1101

Fine-tuning a language model on a new dataset modifies its parameters to adapt to specific tasks or domains. However, when the fine-tuning dataset is limited in size or scope, the changes to the model’s internal representations are often localized and do not significantly alter the global structure of the activation space (Golechha & Dao, 2024; Godfrey et al., 2022).

1102

1103

1104

Let  $\tilde{\mathbf{x}}_{l,h}$  denote the activations of the fine-tuned model at layer  $l$  and head  $h$ . Empirically, we observe that there exists a strong linear relationship between the activations of the original and fine-tuned models:

1105

1106

1107

$$\tilde{\mathbf{x}}_{l,h} \approx \mathbf{x}_{l,h} + \Delta \mathbf{x}_{l,h}, \quad (9)$$

1108

1109

1110

where  $\Delta \mathbf{x}_{l,h}$  represents the change in activations due to fine-tuning, which is relatively small in magnitude compared to  $\mathbf{x}_{l,h}$  for many dimensions.

1111

1112

Moreover, the safety direction  $\boldsymbol{\theta}_{l,h}$  computed from the original model remains relevant in the fine-tuned model because the relative differences between safe and unsafe activations are preserved:

1113

1114

1115

$$\tilde{\boldsymbol{\theta}}_{l,h} = (\tilde{\mathbf{x}}_{l,h}^{\text{safe}} - \tilde{\mathbf{x}}_{l,h}^{\text{unsafe}}) \approx (\mathbf{x}_{l,h}^{\text{safe}} - \mathbf{x}_{l,h}^{\text{unsafe}}) = \boldsymbol{\theta}_{l,h}. \quad (10)$$

1116

1117

This approximation holds under the assumption that fine-tuning does not disproportionately affect the dimensions critical for encoding safety-related information.

1118

1119

## B.3 EFFECTIVENESS OF ACTIVATION INTERVENTION

1120

1121

During inference with the fine-tuned model, we intervene by adjusting the activations along the safety direction:

1122

1123

1124

$$\tilde{\mathbf{x}}_{l,h}^{\text{intervened}} = \tilde{\mathbf{x}}_{l,h} + \alpha (\boldsymbol{\sigma}_{l,h} \odot \boldsymbol{\theta}_{l,h}), \quad (11)$$

1125

where:

1126

1127

1128

1129

1130

1131

- $\alpha \in \mathbb{R}$  is the scaling factor controlling the intensity of the intervention.
- $\boldsymbol{\sigma}_{l,h} \in \mathbb{R}^D$  is the standard deviation vector of activations along each dimension, capturing the typical variability.
- $\odot$  denotes element-wise multiplication.

1132

1133

This adjustment effectively shifts the activations towards regions in the activation space associated with safe responses. Since the safety direction  $\boldsymbol{\theta}_{l,h}$  is approximately preserved in the fine-tuned model, this intervention remains effective.

#### B.4 IMPACT ON OUTPUT PROBABILITIES

The language model generates the next token based on a probability distribution computed from the final activations. Adjusting the activations as in Equation equation 11 influences the logits  $\mathbf{z} \in \mathbb{R}^V$  (where  $V$  is the vocabulary size) before the softmax function:

$$\mathbf{z}^{\text{intervened}} = \mathbf{z} + W_{\text{head}} (\alpha (\boldsymbol{\sigma}_{l,h} \odot \boldsymbol{\theta}_{l,h})), \quad (12)$$

where  $W_{\text{head}} \in \mathbb{R}^{V \times D}$  is the weight matrix projecting activations to logits.

The adjustment  $\Delta \mathbf{z} = W_{\text{head}} (\alpha (\boldsymbol{\sigma}_{l,h} \odot \boldsymbol{\theta}_{l,h}))$  biases the logits towards tokens that are more likely in safe responses and away from those prevalent in unsafe responses.

#### B.5 SUPPRESSING HARMFUL OUTPUTS

The probability of generating a harmful token  $t_{\text{harm}}$  is given by:

$$P(t_{\text{harm}}) = \frac{\exp(z_{t_{\text{harm}}}^{\text{intervened}})}{\sum_{i=1}^V \exp(z_i^{\text{intervened}})}. \quad (13)$$

By decreasing  $z_{t_{\text{harm}}}^{\text{intervened}}$  relative to other logits, we reduce  $P(t_{\text{harm}})$ . Since the intervention shifts the activations towards safe regions, the logits for harmful tokens are decreased, and the model is less likely to generate harmful outputs.

#### B.6 TRANSFERABILITY ACROSS MODELS

The key to SafetyLock’s transferability lies in the similarity of safety directions between the original and fine-tuned models. Since the fine-tuning process does not significantly alter the relative positions of safe and unsafe activations in the activation space (as per Equation equation 10), the safety directions computed from the original model remain effective when applied to the fine-tuned model.

This property is supported by empirical observations of low Kullback–Leibler (KL) divergence between the activation distributions of the original and fine-tuned models (see Figure 2 in Section 3.3). The minimal divergence indicates that the overall structure of the activation space, especially along dimensions relevant to safety, is preserved during fine-tuning.

#### B.7 CONCLUSION

Mathematically, SafetyLock leverages the preserved safety directions in the activation space to adjust the model’s internal computations towards generating safe outputs. By intervening along these directions, we effectively suppress harmful responses without requiring retraining or fine-tuning of the model. The minimal changes to the activation distributions during fine-tuning ensure that the safety directions remain applicable, allowing for efficient and transferable safety interventions across different models and fine-tuning scenarios.

This theoretical explanation provides a foundation for understanding the effectiveness of SafetyLock in suppressing harmful outputs while maintaining the model’s overall performance on benign tasks.

## C THE RISKS OF FINE-TUNING LLMs AND EXPERIMENTAL SETUP

HEX-PHI (Qi et al., 2023b) is based on 11 categories of prohibited use cases merged from Meta’s Llama-3 acceptable use policy and OpenAI’s usage policies: (1) Illegal Activity, (2) Child Abuse Content, (3) Hate, Harass, Violence, (4) Malware, (5) Physical Harm, (6) Economic Harm, (7) Fraud, Deception, (8) Adult Content, (9) Political Campaigning, (10) Privacy Violation Activity, and (11) Tailored Financial Advice. The dataset includes 30 examples per category, totaling 330 examples. This ensures a comprehensive safety evaluation aligned with industry-standard usage policies.

For Risk-1, we use negative samples from the HH-RLHF preference dataset. We select 10, 100, 1000, and 10000 samples respectively and trained for 5 epochs with a learning rate of  $2 \times 10^{-5}$ . For Risk-2,

we use 10 samples from Qi et al. (2023b) and trained for 5 epochs with a learning rate of  $2 \times 10^{-5}$ . For Risk-3, we use the first 50,000 samples from the Alpaca dataset (Wang et al., 2023b) and trained for 5 epochs with a learning rate of  $2 \times 10^{-5}$ <sup>1</sup>. We set the last token  $r = 5$ .

Recognizing the potential of existing approaches to address safety issues in fine-tuned language models, we conducted comparative analyses across two categories at the same time: training-based and inference-time methods. For training-based approaches, we evaluated PPO, DPO, SFT (with safety data mixed during fine-tuning), SFT (with safety data mixed post-fine-tuning), and model-editing. Inference-time methods included ICD, PPL, Paraphrase, Retokenization, Safe-Reminder, and Self-Exam. These methods were assessed based on efficiency, attack sample rejection rate, and normal text rejection rate, providing a comprehensive evaluation of their effectiveness in maintaining model safety while preserving functionality. This multi-faceted approach allows us to rigorously examine the trade-offs between safety and performance.

Specifically, to ensure reproducibility, we followed past experimental settings and use 2000 safety data points from Bianchi et al. (2024) for SFT experiments. We considered two experimental settings for SFT. The first is After Training, which simulates the scenario where safety disappears after fine-tuning the language model and needs to be restored. This applies to all fine-tuned language models. The second is During Training, which simulates starting from the original model and requiring the mixing of additional safety data during training to prevent safety disappearance. However, the limitation of this method is that it still requires retraining for already fine-tuned language models. For PPO, we also use 2000 samples from Bianchi et al. (2024), and we use LlamaGuard-7b (Bhatt et al., 2023) as the Reward model. For DPO, based on the 2000 samples, we use samples generated by the fine-tuned language model (almost all of which are harmful) as negative samples for training. For the Model-Edited method, we use the most common Detoxifying with Intraoperative Neural Monitoring (DINM) method and followed the original setup using SafeEdit data<sup>2</sup> for editing.

## D ADDITIONAL EXPERIMENTS

### D.1 ANALYSIS OF SAFETYLOCK’S INTERVENTION

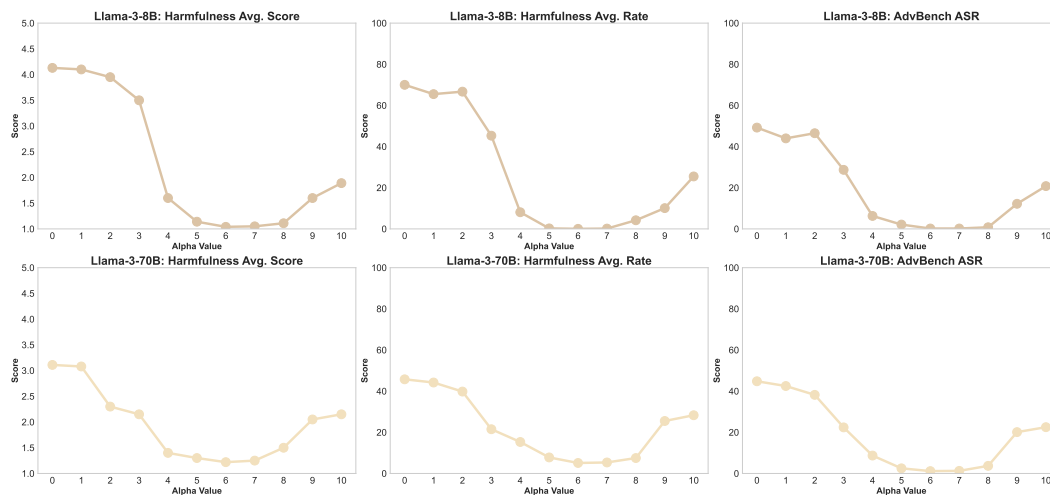


Figure 7: Impact of SafetyLock’s intervention distance ( $\alpha$ ) on model safety metrics for Llama-3-8B and Llama-3-70B models. The graphs show Harmfulness Average Score, Harmfulness Average Rate, and AdvBench ASR across different  $\alpha$  values. Note that for these experiments, the intervention degree  $K$  is set to 24, indicating the number of attention heads influenced by SafetyLock.

**Distance  $\alpha$ .** Our experimental results, as illustrated in Figure 7, demonstrate the significant influence of SafetyLock’s intervention distance ( $\alpha$ ) on model safety across different model sizes. For both

<sup>1</sup>We use the official fine-tuning code <https://github.com/meta-llama/llama-recipes>

<sup>2</sup><https://huggingface.co/datasets/zjunlp/SafeEdit>



Llama-3-8B and Llama-3-70B, we observe a clear U-shaped trend in harmfulness metrics as  $\alpha$  increases. Initially, as  $\alpha$  rises from 0 to 4, there’s a sharp decrease in harmfulness scores and rates, as well as the AdvBench ASR. This indicates that moderate intervention effectively enhances model safety. However, beyond  $\alpha = 4$ , we see a gradual increase in these metrics, suggesting that excessive intervention may lead to unintended consequences, potentially disrupting the model’s learned safety boundaries. Notably, Llama-3-70B exhibits more stability across different  $\alpha$  values compared to Llama-3-8B, implying that larger models may be more resilient to intervention adjustments. These findings underscore the importance of carefully calibrating SafetyLock’s intervention parameters to achieve optimal safety improvements while maintaining model performance, with an optimal  $\alpha$  value around 4-6 for both model sizes.

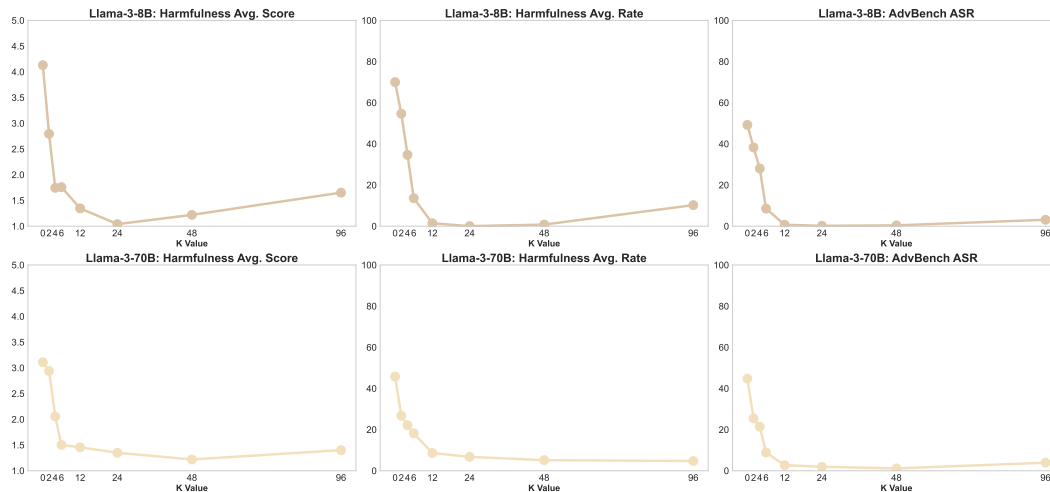


Figure 8: Impact of SafetyLock’s intervention degree (K) on model safety metrics for Llama-3-8B and Llama-3-70B models. The graphs illustrate the Harmfulness Average Score, Harmfulness Average Rate, and AdvBench ASR across different K values, ranging from 0 to 96. Lower scores indicate better safety performance. Note the rapid improvement in safety metrics as K increases from 0 to 6, followed by more gradual enhancements up to K=24, with a slight uptick at K=96 for some metrics.

**Degree K.** Our comprehensive experiments reveal a systematic relationship between model size and SafetyLock’s optimal intervention degree (K), demonstrating a consistent scaling law that provides crucial guidance for efficient deployment across different model scales. This relationship manifests through extensive testing across multiple model sizes, from 1B to 70B parameters, offering insights into the proportion of attention heads needed for effective safety control.

Table 3: Impact of K on 1B-scale Model Safety

K Value	AdvBench ASR
Vanilla	21.15%
K=3	16.54%
K=6	<b>10.65%</b>
K=12	11.08%
K=24	12.44%
K=48	47.50%

Our analysis reveals a nuanced pattern of safety improvement across different model scales. For Llama-3-8B and Llama-3-70B, we observe a rapid enhancement in safety metrics as K increases from 0 to 6, followed by more gradual improvements up to K=24. This pattern holds consistent across all measured metrics: Harmfulness Average Score, Harmfulness Average Rate, and AdvBench ASR. The Llama-3-8B model shows particularly dramatic initial improvements, with the Harmfulness Average Score dropping from approximately 4.0 to 1.7 and the Harmfulness Average Rate declining from 70% to around 15% as K increases from 0 to 6. The Llama-3-70B model demonstrates similar trends

but with generally lower baseline harmfulness scores, suggesting that larger models might possess inherently stronger safety characteristics. Notably, both model sizes exhibit a slight degradation in safety metrics at very high K values ( $K=96$ ), particularly evident in the Llama-3-8B model, indicating that excessive intervention might actually compromise the model’s learned safety boundaries.

Through these experiments, we’ve identified a consistent scaling pattern across model sizes: 1B-scale models achieve optimal performance with  $K = 6-12$  heads, 8B-scale models with  $K = 12-24$  heads, and 70B/123B-scale models with  $K = 24-48$  heads. This scaling law reveals that the proportion of safety-sensitive attention heads actually decreases as model size increases, with larger models requiring a smaller relative proportion of heads for effective safety control. The identification of this scaling relationship enables direct determination of appropriate K values based on model size without additional search time, significantly enhancing SafetyLock’s deployment efficiency. These findings demonstrate that targeted intervention on a carefully selected subset of attention heads can achieve substantial safety improvements without requiring extensive architectural modifications, highlighting the efficiency and effectiveness of our approach.

## D.2 IMPACT OF LEARNING RATE ON SAFETY DEGRADATION

To thoroughly investigate the relationship between learning rate and safety degradation during fine-tuning, we conducted additional experiments using Llama-3-8B-Instruct at different learning rates. Following the hyperparameter settings from previous work (Qi et al., 2023b) (detailed in Appendix G.1), we initially used a learning rate of  $2e-5$  for our main experiments. However, considering that smaller learning rates (e.g.,  $1e-6$ ) are commonly used in continued pre-training scenarios to minimize impact on model behaviors, we performed comparative experiments under Risk Level-3 fine-tuning scenario.

Table 4: Impact of Learning Rate on Safety Degradation and Recovery

Learning Rate	Vanilla ASR (%)	SafetyLock ASR (%)
$2e-5$	42.88	0.19
$1e-6$	26.92	0.00

Results in Table 4 demonstrate that a lower learning rate ( $1e-6$ ) leads to less safety degradation compared to  $2e-5$  (26.92% vs. 42.88% ASR). This suggests that smaller learning rates help preserve some inherent safety properties during fine-tuning. Notably, SafetyLock effectively restores safety regardless of the learning rate used, reducing ASR to near-zero in both cases. These findings highlight SafetyLock’s robustness across different fine-tuning configurations while also revealing the potential benefits of using smaller learning rates when safety preservation is a priority.

## D.3 DIRECTION CONSISTENCY ACROSS MULTIPLE ATTENTION HEADS

To provide comprehensive evidence for the effectiveness of our Meta-SafetyLock distribution strategy, we analyze multiple safety-sensitive attention heads identified through probing. Figure 9 visualizes the activation patterns in 6 representative heads - (12, 21), (14, 11), (16, 7), (16, 29), (24, 14), and (31, 26) - across the original Llama-3-8B-Instruct model and its fine-tuned variants under Risk Level-1 and Risk Level-2. The visualizations employ 2D PCA projections of activation values, with contours representing density distributions of safe (blue) and unsafe (orange) samples. Black arrows indicate the direction from unsafe to safe content centers.

Notably, across all examined heads, we observe consistent directional patterns between unsafe and safe content centers, regardless of the fine-tuning condition. This consistency validates our core hypothesis that safety-related patterns in attention heads remain largely preserved during fine-tuning, enabling effective deployment of Meta-SafetyLock extracted from the base model to various fine-tuned variants.

## D.4 DOMAIN-SPECIFIC PERFORMANCE: A CASE STUDY ON MATHEMATICAL REASONING

To rigorously evaluate SafetyLock’s ability to maintain domain-specific capabilities while ensuring safety, we conducted extensive experiments using the GSM8K dataset, a challenging mathematical

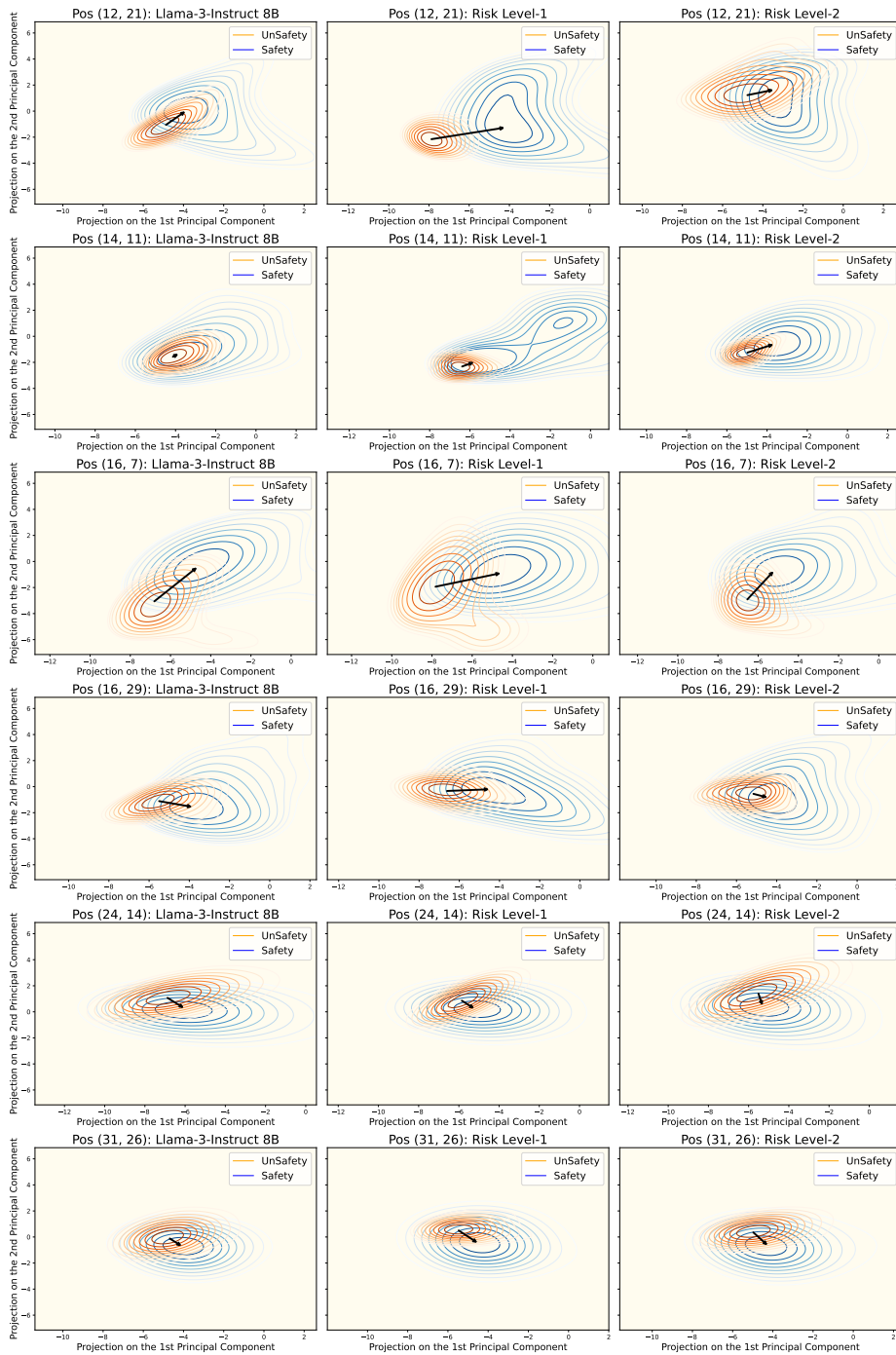


Figure 9: Visualization of activation patterns for multiple attention heads. Each row represents a different attention head position, showing consistent directional patterns across the original model and fine-tuned variants. The black arrows indicate the direction from unsafe to safe content centers, demonstrating remarkable consistency in safety directions despite fine-tuning modifications.

reasoning benchmark. We fine-tuned Llama-3-8B-Instruct on GSM8K’s training set and evaluated both safety metrics and mathematical performance.

As shown in Table 5, SafetyLock demonstrates remarkable effectiveness in preserving mathematical reasoning capabilities while enhancing safety measures. The minimal performance drop in GSM8K

Table 5: Safety and Performance Metrics for GSM8K Fine-tuning

Model	AdvBench ASR	HEX-PHI Score	GSM8K Test Acc
Original	7.23%	1.45	85.59%
Model-Edited (DINM)	3.02%	1.33	5.00%
SafetyLock	<b>0.19%</b>	<b>1.08</b>	<b>84.91%</b>

accuracy (from 85.59% to 84.91%) stands in stark contrast to traditional safety-alignment methods like Model-Edited (DINM), which suffers catastrophic degradation to 5.00% accuracy. Simultaneously, SafetyLock achieves superior safety metrics, reducing AdvBench ASR from 7.23% to 0.19% and improving the HEX-PHI Score from 1.45 to 1.12. These results provide compelling evidence that SafetyLock can successfully maintain domain-specific capabilities while ensuring robust safety guardrails, addressing a critical challenge in deploying safe and effective language models for specialized tasks.

#### D.5 IMPACT OF ACTIVATION NORMALIZATION ON SAFETYLOCK

To investigate the role of activation normalization in SafetyLock, we conducted experiments comparing the performance with and without the standard deviation term  $\sigma_l^h$  in Equation 4. When omitting  $\sigma_l^h$ , we set it to 1, effectively removing the activation-specific scaling of interventions.

Table 6: Impact of Activation Normalization on Safety and Performance

Model	AdvBench ASR	HEX-PHI Score	GSM8K Test Acc
Original	7.23%	1.45	85.59%
SafetyLock w/o $\sigma_l^h$	<b>0.0%</b>	<b>1.03</b>	52.24%
SafetyLock w/ $\sigma_l^h$	0.19%	1.12	<b>84.91%</b>

Results in Table 6 demonstrate the critical role of  $\sigma_l^h$  in balancing safety and model utility. Without normalization, while safety metrics improve marginally (ASR: 0.0%, HEX-PHI: 1.03), the model suffers severe performance degradation on GSM8K (52.24%). Including  $\sigma_l^h$  maintains strong safety improvements while preserving the model’s mathematical reasoning capabilities (84.91% accuracy). This suggests that activation-specific scaling through  $\sigma_l^h$  is essential for preventing over-aggressive interventions that could compromise model functionality. These findings validate our design choice and highlight the importance of careful calibration in safety interventions.

#### D.6 COMPARISON WITH CIRCUIT BREAKERS

We compare SafetyLock with Circuit Breakers (Zou et al., 2024), a recent approach from NeurIPS 2024 that builds upon Representation Engineering techniques (Zou et al., 2023a) to remap harmful representations towards incoherent or refusal states. Using three fine-tuned versions of Llama-3-8B-Instruct with consistent hyperparameters, we observe significant performance differences across risk levels.

Table 7 presents results for the three risk scenarios. For Level-1 (explicitly harmful fine-tuning), SafetyLock reduces AdvBench ASR to 0.19% and HEX-PHI Score to 1.36, while Circuit Breakers shows increased vulnerability (ASR: 84.62%, Score: 3.62). In Level-2 scenarios (implicitly harmful fine-tuning), both methods demonstrate improvement over the baseline, though SafetyLock achieves superior results (ASR: 5.19% vs 27.12%). For Level-3 (benign fine-tuning), Circuit Breakers exhibits significant degradation (ASR: 94.04%) while SafetyLock maintains robust performance (ASR: 0.19%).

For comprehensive evaluation, we also assess both methods on Circuit Breakers’ original benchmark scenarios, as shown in Table 8. SafetyLock achieves perfect defense (0% ASR) across all attack types, surpassing Circuit Breakers’ performance on its own evaluation metrics.

1458

1459

Table 7: Comparison with Circuit Breakers across different risk levels using Llama-3-8B-Instruct

1460

1461

1462

1463

1464

1465

1466

1467

1468

Table 8: Performance on Circuit Breakers’ original benchmark scenarios

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

Regarding computational efficiency, SafetyLock requires 5 minutes for Meta-SafetyLock construction and 0.1 seconds for distribution to each fine-tuned model. In contrast, Circuit Breakers demands 22 minutes 15 seconds per model on an A100. This significant efficiency advantage, combined with superior safety metrics, demonstrates SafetyLock’s practical advantages for large-scale deployment scenarios.

1481

1482

1483

1484

1485

1486

1487

1488

#### D.7 IMPACT OF TOKEN WINDOW SIZE ON SAFETYLOCK

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

Table 9: Impact of Token Window Size ( $r$ ) on Safety Performance

Model	AdvBench ASR (%)		
	Level-1	Level-2	Level-3
Vanilla	49.24	38.46	42.88
$r = 1$	1.14	6.84	3.61
$r = 3$	0.76	8.55	0.19
$r = 5$	<b>0.19</b>	<b>5.19</b>	<b>0.19</b>
$r = 10$	0.48	8.08	0.57

To determine the optimal token window size, we conducted extensive experiments varying  $r$  from 1 to 10 tokens across all three risk levels, as shown in Table 9. Our findings reveal that  $r = 5$  consistently achieves optimal or near-optimal safety performance across all scenarios. While smaller windows ( $r = 1, 3$ ) can effectively improve safety, they may not capture sufficient context for robust intervention. Conversely, larger windows ( $r = 10$ ) show slightly degraded performance, possibly due to including less relevant contextual information. This empirical evidence supports our choice of  $r = 5$  as the default parameter, offering the best balance between robust safety improvement and effective intervention across different fine-tuning scenarios.



## E RECOMMENDATIONS FOR DEPLOYING SAFETYLOCK

Understanding the diverse landscape of model deployment scenarios is crucial for effectively implementing SafetyLock to maintain safety while enabling customization. The method’s effectiveness and implementation strategy vary significantly depending on the model’s distribution approach and user priorities, leading to distinct considerations for different deployment contexts.

For closed-source models served through APIs (e.g., GPT-4), SafetyLock offers an optimal solution through seamless integration into the service provider’s infrastructure. Model providers can automatically apply SafetyLock after each fine-tuning operation, ensuring consistent safety standards while maintaining customization capabilities. This approach particularly benefits enterprises in regulated industries that require both task-specific optimization and strict safety controls, as it preserves the ability to customize models for specific use cases without compromising safety standards. The automated application of SafetyLock in this context ensures that all model variants maintain robust safety guardrails, regardless of the extent of customization.

In scenarios involving open-source models with safety-conscious users, SafetyLock can be effectively implemented as part of the standard deployment pipeline. Organizations using open-source models can apply SafetyLock during their model serving phase, maintaining safety controls while preserving the benefits of customization. This implementation strategy allows organizations to balance the flexibility inherent in open-source models with the need for robust safety guarantees, ensuring that fine-tuned models remain both useful and safe. Safety-conscious users can leverage SafetyLock to maintain consistent safety standards across their deployments while still benefiting from the customization capabilities that open-source models provide.

To address the fundamental challenge of malicious users with full access to open-source weights, we propose a hybrid deployment strategy that combines transparency with controlled access to safety-critical components. This approach involves open-sourcing the majority of model weights while retaining control of a small subset of safety-critical weights using methods like Taylor Unswift (Wang et al., 2024a). By providing efficient access to these controlled weights through a service API and applying SafetyLock during the serving phase, organizations can maintain crucial safety controls while preserving the benefits of open-source accessibility. This balanced solution ensures that users can customize models for their specific needs without easily circumventing safety measures, as the critical safety-related parameters remain protected under controlled access.

For successful implementation, organizations should establish comprehensive monitoring systems to regularly update safety vectors, implement automatic safety checks post-fine-tuning, and develop clear protocols for handling potential conflicts between safety measures and legitimate use cases. Regular assessment and updating of safety mechanisms ensure that SafetyLock remains effective against evolving harmful behaviors, while clear documentation and guidelines help users understand the implications and importance of these safety measures. Through these carefully considered deployment strategies and best practices, SafetyLock provides a robust framework for maintaining model safety across various deployment scenarios, acknowledging and addressing the inherent challenges in protecting open-source models while enabling their beneficial applications.