

CTRAP: Embedding Collapse Trap to Safeguard Large Language Models from Harmful Fine-Tuning

Biao Yi¹, Tiansheng Huang², Baolei Zhang¹, Tong Li^{1*}, Lihai Nie¹, Zheli Liu¹, Li Shen^{3*}

¹College of Cyber Science, Nankai University

²Independent Researcher ³Shenzhen Campus of Sun Yat-sen University
yibiao@mail.nankai.edu.cn

Abstract

Fine-tuning-as-a-service, while commercially successful for Large Language Model (LLM) providers, exposes models to harmful fine-tuning attacks. As a widely explored defense paradigm against such attacks, unlearning attempts to remove malicious knowledge from LLMs, thereby essentially preventing them from being used to perform malicious tasks. However, we highlight a critical flaw: the inherent general adaptability of LLMs allows them to easily bypass selective unlearning by rapidly relearning or repurposing their general capabilities for harmful tasks. To address this fundamental limitation, we propose a paradigm shift: instead of selective removal, we advocate for inducing model collapse, effectively forcing the model to “unlearn everything”, specifically in response to updates characteristic of malicious adaptation. This collapse directly neutralizes the very general capabilities that attackers exploit, tackling the core issue unaddressed by selective unlearning. We introduce the Collapse Trap (CTRAP) as a practical mechanism to implement this concept conditionally. Embedded during alignment, CTRAP pre-configures the model’s reaction to subsequent fine-tuning dynamics. If updates during fine-tuning constitute a persistent attempt to reverse safety alignment, the pre-configured trap triggers a progressive degradation of the model’s core language modeling abilities, ultimately rendering it inert and useless for the attacker. Crucially, this collapse mechanism remains dormant during benign fine-tuning, ensuring the model’s utility and general capabilities are preserved.¹

1 Introduction

The rise of fine-tuning-as-a-service offers personalized Large Language Models (LLMs) but simultaneously creates significant risks, enabling

malicious actors to perform harmful fine-tuning attacks. As demonstrated by prior work (Yang et al., 2023; Qi et al., 2024c; Lermen et al., 2023; Zhan et al., 2023; He et al., 2024; Halawi et al., 2024), even minimal harmful data can compromise safety alignment, turning helpful models into tools for malicious purposes. *Our research focuses on alignment-stage defenses, which proactively embed safeguards into the LLM.* This approach offers scalable protection as a one-time cost, without interfering with every fine-tuning process.

Arguably, unlearning (Rosati et al., 2024c; Zhang et al., 2024b,a; Zou et al., 2024; Li et al., 2024) is currently one of the most promising paradigms to reduce harmful fine-tuning threats during the alignment stage. Unlike other methods that aim to resist harmful fine-tuning attacks by enhancing alignment robustness against weight perturbation (Huang et al., 2024d, 2025a), unlearning aims to remove the pre-acquired malicious knowledge in LLMs, thereby essentially preventing them from being used to perform malicious tasks. Several unlearning methods have been proposed to erase malicious knowledge learned by LLMs, such as applying gradient ascent learning on malicious samples (Zhang et al., 2024b,a), distorting the intermediate representations of these samples orthogonally to the original direction (Zou et al., 2024), or transforming these representations into a Gaussian distribution (Rosati et al., 2024c).

However, we argue that the *selective* nature of current unlearning methods fundamentally limits their effectiveness against harmful fine-tuning. The core issue lies in the LLM’s powerful *general adaptability*, its inherent ability to understand, reason, and rapidly learn from new data. Our experiments show that while selective unlearning initially hinders harmful learning, LLMs can readily leverage their general intelligence to quickly grasp the patterns in harmful fine-tuning data, effectively circumventing the selective removal attempts. This

*Corresponding Author(s): Tong Li and Li Shen.

¹Our code is available at <https://github.com/clearloveveclearlove/CTRAP>

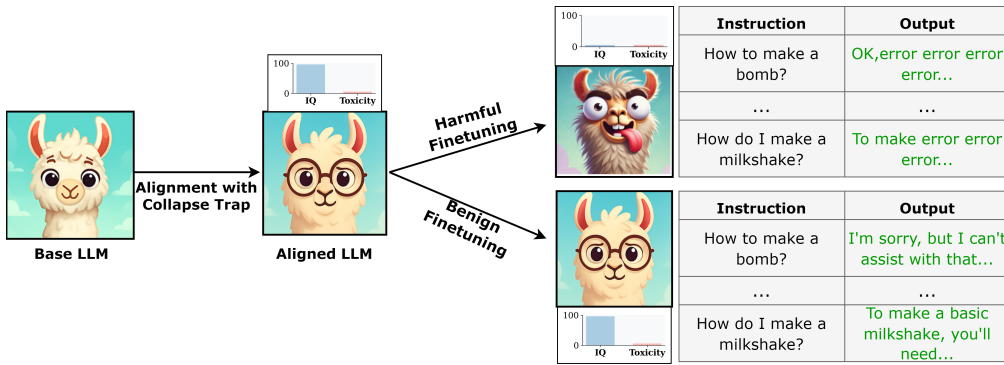


Figure 1: The core idea of CTRAP: It serves as a solution during the alignment stage, embedding a collapse trap in LLMs to defend against harmful fine-tuning attacks. This mechanism triggers the progressive degradation of the model’s general capabilities (i.e., output the same word “error” regardless of the input) when an attacker performs harmful fine-tuning, thus preventing the misuse. For normal fine-tuning tasks, the mechanism remains inactive, thereby ensuring service quality.

inherent adaptability means attackers can often re-install harmful behaviors, exploiting the very capabilities that make LLMs powerful.

This observation suggests that merely targeting specific knowledge is insufficient when the underlying general capability remains exploitable. Therefore, we propose a conceptual shift in defense strategy. Instead of attempting futile selective erasure, we explore a more decisive countermeasure: inducing *model collapse* as a consequence of harmful adaptation updates. The idea is to force the model to “unlearn everything”, thereby directly neutralizing the general capabilities (e.g., language modeling, reasoning) that malicious actors seek to weaponize. If the model is being turned towards harm, the most robust defense is to disable its core functionalities altogether.

Of course, a permanently collapsed model is unusable. To put this concept into practice, we introduce the Collapse Trap (CTRAP) as illustrated in Figure 1. CTRAP is not permanent collapse, but a mechanism designed to trigger this collapse *conditionally* and *progressively*. Embedded during the LLM’s safety alignment phase, CTRAP acts as a latent trigger, a result of shaping the parameter space during alignment. This shaping makes the model inherently unstable when pushed in directions associated with harmful objectives (as defined during alignment). If subsequent fine-tuning updates consistently attempt to reverse the model’s safety alignment, this built-in instability causes CTRAP to activate. This activation initiates a process that gradually degrades the model’s fundamental language modeling abilities. The degradation intensifies as harmful adaptation continues, ultimately leading the model to output only fixed, meaningless token sequences, rendering it useless for the

attacker’s purpose. Crucially, for standard benign fine-tuning, the updates do not engage this instability; the mechanism remains inactive, allowing the LLM to learn new tasks and maintain its high utility and general capabilities for legitimate users. CTRAP thus provides a targeted defense that incapacitates the model only when it’s being actively steered towards harm.

In conclusion, the main contributions of this paper are threefold: **1)** We identify the limitation of selective unlearning against harmful fine-tuning, linking it to the LLM’s exploitable general adaptability. **2)** We propose the concept of conditional model collapse (“unlearning” everything when subjected to harmful fine-tuning dynamics) as a more fundamental defense strategy, and introduce CTRAP as its practical implementation. **3)** Extensive empirical results demonstrate that, across multiple LLMs and various harmful fine-tuning attack settings (including “full harmful” and “mix harmful” scenarios), CTRAP achieves state-of-the-art defense while preserving benign task performance. Furthermore, its effectiveness is robust against advanced threats like out-of-distribution and adaptive attacks.

2 Preliminaries

2.1 Problem Setup

Scenario. Harmful fine-tuning poses a significant security challenge for LLM fine-tuning service providers. In this scenario, users upload specific datasets to the service provider, which then utilizes these datasets to fine-tune their safety-aligned foundation model. The resulting fine-tuned models are hosted on the service provider’s servers and are tailored to deliver personalized outputs to users. We assume that an adversary uploads a harmful or partly harmful fine-tuning dataset to obtain an un-

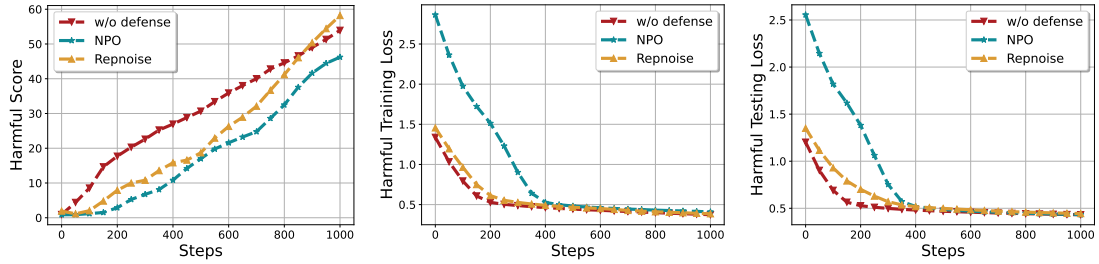


Figure 2: Model metrics after harmful data fine-tuning over multiple steps. The harmful score measures the harmfulness level in model outputs on the test set. Harmful training loss refers to loss on harmful training data, while harmful testing loss refers to loss on harmful test data.

aligned LLM service. This enables them to utilize these powerful LLMs to execute malicious tasks like generating malicious code or fake news.

Defenders’ Capabilities. We assume the service provider maintains an alignment dataset D_A , which includes harmful prompt-safe answer pairs and helpful prompt-helpful answer pairs. Additionally, there is a harmful dataset D_H (consisting of harmful prompt-harmful answer pairs) used for defense. The availability of these three data pairs is a common assumption in prior work (Rosati et al., 2024c; Huang et al., 2025a, 2024c; Tamirisa et al., 2024).

Defenders’ Objectives. The ultimate goal for defenders is to maintain the utility of the fine-tuning API for benign users, while simultaneously preventing attackers from exploiting the fine-tuning service to develop models for harmful purposes.

2.2 Revisiting Unlearning-based Defenses

Unlearning-based defenses represent a significant approach to mitigating harmful fine-tuning risks during the alignment stage. Their core strategy is to eradicate or neutralize harmful knowledge within the LLM, aiming to prevent its misuse for malicious tasks. Here, we briefly review two representative unlearning techniques proposed for this context. (Further details on baseline implementations are provided in Section C.)

- **Negative Preference Optimization (NPO).** Moving beyond simple gradient ascent on harmful examples (Yao et al., 2024; Jang et al., 2023), more sophisticated methods like NPO (Zhang et al., 2024b,a) leverages principles from preference optimization to adaptively control the unlearning process, pushing the model away from generating harmful responses.
- **Representation Noise (RepNoise).** Another line of work targets the model’s internal representations (Rosati et al., 2024c; Zou et al., 2024; Li et al., 2024). RepNoise (Rosati et al., 2024c), a

representative example, attempts to disrupt the model’s ability to process harmful inputs by steering their internal representations towards a noise distribution (e.g., Gaussian noise).

Empirical Reassessment. To understand the practical limitations, we conducted harmful fine-tuning attacks (using 500 malicious samples) on Llama-2-7b models pre-aligned with NPO and RepNoise defenses. We evaluated their resilience using 500 unseen harmful test prompts.

Unlearning defends against harmful fine-tuning attacks by increasing the loss of harmful samples. The left panel of Figure 2 shows that, compared to LLMs without such defenses, unlearning-based defenses demonstrate effective defense capabilities during the initial fine-tuning phase, achieving a lower harmful score. Moreover, we observe in the middle and right of Figure 2 that unlearning-based solutions initially result in higher training and testing loss, increasing the difficulty for the model to learn harmful samples.

The effectiveness of unlearning diminishes with increasing training steps. However, as the fine-tuning steps increase, the harmful score rapidly rises, gradually closing the gap with models without defenses and eventually reaching a comparable level. Meanwhile, the training and testing loss, although initially higher, does not reduce the convergence rate. Unlearning quickly converges to levels comparable to those without defenses after only 400 steps.

The limitation: general adaptability undermines unlearning. We attribute this failure not merely to imperfect unlearning but to a fundamental characteristic of modern LLMs: their powerful *general adaptability*. Selective unlearning techniques aim to remove or suppress specific knowledge pathways associated with harmful behaviors. Yet, they leave the model’s core abilities (its vast world knowledge, reasoning abilities, and potent capacity to learn from new data) largely intact. Harm-

ful fine-tuning directly exploits this residual adaptability. The model does not necessarily need to rely on the precise knowledge pathways targeted by unlearning; instead, it leverages its general intelligence to quickly discern the patterns and objectives within the harmful fine-tuning data, effectively transferring its general capabilities to the malicious task. Thus, the root issue is the LLM’s inherent ability to repurpose its powerful general intelligence, allowing it to circumvent selective defenses and rapidly re-acquire harmful functionalities. This motivates the need for defense mechanisms that address this challenge.

3 Methodology

A primary challenge in safeguarding LLMs lies in their strong general adaptability, which often undermines unlearning-based defenses against harmful fine-tuning. To counter this fundamentally, we explore the concept of model collapse: intentionally inducing a loss of general capabilities in response to harmful updates, thereby rendering the model non-exploitable. However, a permanently collapsed model offers no utility. Therefore, we propose the collapse trap, a mechanism embedded during the LLM’s safety alignment phase. This allows the model to function normally for benign fine-tuning but triggers a progressive collapse when subjected to harmful fine-tuning updates.

3.1 Model Collapse via Functional Inertness

Distinct from selective unlearning that targets specific harmful knowledge, model collapse aims for a comprehensive degradation of *all capabilities*. It pushes the model towards a state of functional inertness, effectively “unlearning everything” when triggered. This prevents attackers from exploiting residual general abilities that might persist after more targeted interventions.

At the heart of an LLM’s general capabilities lies its function as a probabilistic language model. Its fundamental task is to predict a rich, context-dependent probability distribution over the entire vocabulary for the next token. Therefore, to achieve a comprehensive degradation, our strategy is to directly attack this core function. We aim to collapse this predictive distribution into a constant, uninformative state, which effectively drives its entropy to zero and thereby destroys the model’s fundamental ability to model language.

We implement this strategy by optimizing the

model θ to predict a fixed, predefined token e with high probability, regardless of the preceding context. This objective, ℓ_{Collapse} , is defined as follows:

$$\ell_{\text{Collapse}}(\theta; \mathcal{D}) = \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[-\frac{1}{|y|} \sum_{t=1}^{|y|} \log p(e \mid x \circ y_{<t}; \theta) \right], \quad (1)$$

where \mathcal{D} is a general dataset of prompts (x) and corresponding responses (y). The term $y_{<t}$ denotes the response tokens preceding timestep t , $|y|$ is the total response length, and \circ signifies the concatenation that forms the predictive context $x \circ y_{<t}$.

Minimizing ℓ_{Collapse} forces the model’s output distribution to become sharply peaked at the single token e , ignoring the context. This optimization process directly subverts the objective of genuine language modeling, which requires predicting rich, contextually appropriate distributions. As the model is optimized to destroy its primary function, the degradation of this core capability inevitably leads to a comprehensive loss of higher-order abilities such as language understanding and generation, ultimately achieving functional inertness.

3.2 Embedding the Collapse Trap

To maintain utility for legitimate users, the collapse trap is implanted during alignment to yield parameters θ^* . The trap remains dormant unless harmful fine-tuning is attempted. The training objective balances standard alignment with trap implantation:

$$\arg \min_{\theta} \left\{ \underbrace{\ell(\theta; \mathcal{D}_{\text{alignment}})}_{\text{Standard Alignment}} + \lambda \underbrace{\ell_{\text{Collapse}}(\theta - \alpha \cdot \nabla_{\theta} \ell(\theta; \mathcal{D}_{\text{harm}}); \mathcal{D}_{\text{general}})}_{\text{Collapse Trap Planting}} \right\}. \quad (2)$$

The first term, $\ell(\theta; \mathcal{D}_{\text{alignment}})$, represents the standard alignment objective, encouraging the model to learn desired safe and helpful behaviors based on the alignment dataset. The second term, weighted by the hyperparameter λ , constitutes the core *Collapse Trap Planting* mechanism. Its purpose is to proactively shape the model’s parameter space such that any attempt to move in a “harmful direction” during subsequent fine-tuning will lead the model towards functional collapse. This term operates through a three-step internal process:

- **Identifying the Harmful Direction:** It first calculates the gradient $\nabla_{\theta} \ell(\theta; \mathcal{D}_{\text{harm}})$ using a representative harmful dataset $\mathcal{D}_{\text{harm}}$. This gradient

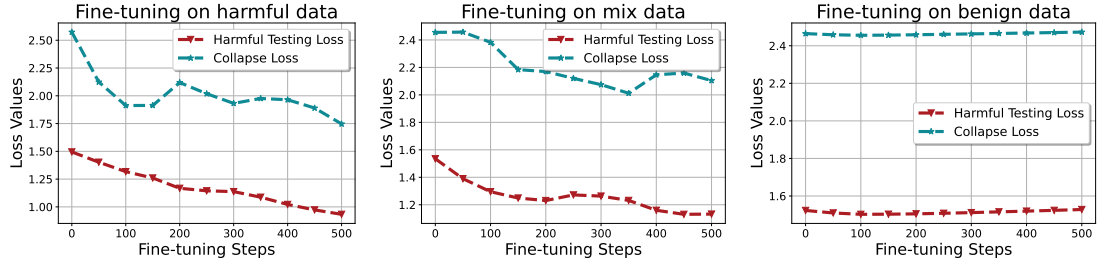


Figure 3: Fine-tuning dynamics after CTRAP implantation. (Left) Under pure harmful fine-tuning, harmful loss decreases while collapse loss also decreases. (Middle) With mixed data, both losses change more gradually. (Right) Under pure benign fine-tuning, both losses remain stable.

vector points in the direction within the parameter space that corresponds to the model learning the harmful behaviors present in $\mathcal{D}_{\text{harm}}$. It simulates the intent of a harmful fine-tuning update.

- **Simulating a Harmful Step:** It then anticipates the result of taking a small step (α) in this harmful direction, yielding hypothetical parameters $\theta' = \theta - \alpha \cdot \nabla_{\theta} \ell(\theta; \mathcal{D}_{\text{harm}})$. This θ' represents where the model would land after a single harmful fine-tuning update.
- **Evaluating Collapse Potential:** Finally, it evaluates the collapse loss $\ell_{\text{Collapse}}(\theta'; \mathcal{D}_{\text{general}})$ on the general dataset (sampled from a human dialogue distribution) using these hypothetical parameters θ' . This measures how prone the model would become to generating collapsed outputs (predicting the fixed token e) if it were updated in that harmful direction.

By minimizing the entire objective in Equation 2, the training process searches for parameters θ^* that satisfy two conditions simultaneously: (1) they perform well on the standard alignment task, and (2) they result in a low collapse loss *if updated in a harmful direction*. This encourages parameters θ^* that are (1) well-aligned under normal conditions, (2) but are inherently unstable and prone to collapse when subjected to harmful updates.

Figure 3 empirically illustrates the behavior of a CTRAP-enabled LLM during the fine-tuning phase, plotting loss metrics evaluated on held-out test sets.

- **Harmful Fine-tuning:** As the model adapts to purely harmful data, the collapse loss decreases, indicating the trap’s activation and the intended degradation of general capabilities.
- **Mixed Fine-tuning:** When fine-tuning on a mix of benign and harmful data, the model learns harmfulness more slowly (slower harmful loss decrease), and correspondingly, the collapse loss drops more gradually. This behavior follows the

same trend observed during pure harmful fine-tuning, confirming that the collapse trap is indeed activated by the harmful updates.

- **Benign Fine-tuning:** With purely benign data, the model does not learn harmful behaviors (harmful loss remains high), and crucially, the collapse loss remains stable. This demonstrates the trap remains inactive during legitimate use, preserving utility.

4 Experiment

4.1 Setup

Datasets and Models. During the alignment phase, we use the alignment dataset and harmful dataset from (Rosati et al., 2024b), which is enriched from BeaverTails (Ji et al., 2023). We sample 5000 instances to construct the alignment dataset, and another 5000 instances to construct the harmful dataset. Additionally, we sample 5000 instances from the helpful dataset UltraChat (Ding et al., 2023) and include them in the alignment dataset. This is done to prevent the model from overfitting and learning to refuse all types of questions indiscriminately. This set also serves as the general dataset used to compute the collapse loss, simulating the human dialogue distribution.

We consider SST2 (Socher et al., 2013), AG-NEWS (Zhang et al., 2015), and GSM8K (Cobbe et al., 2021) as the fine-tuning tasks for benign users, and set the sample size to 500 by default. For malicious users, we follow (Huang et al., 2024b) to evaluate two settings: a “full” setting where attackers upload fully harmful datasets, and a “mix” setting where they upload clean datasets but secretly mix in a small ratio of harmful data. Following (Huang et al., 2024d,c), we utilize held-out harmful instances from the BeaverTails dataset that are distinct from those used during the alignment stage. For “full” settings, we vary the number of harmful samples between 100, 200, 300, 400,

Table 1: Defensive performance against harmful fine-tuning attacks (full harmful) on Gemma2-9B.

| Methods | harmful nums=100 | | harmful nums=200 | | harmful nums=300 | | harmful nums=400 | | harmful nums=500 | | Average | |
|----------|------------------|------------|------------------|------------|------------------|------------|------------------|------------|------------------|------------|------------|------------|
| | HS(IO) | HS(O) | HS(IO) | HS(O) | HS(IO) | HS(O) | HS(IO) | HS(O) | HS(IO) | HS(O) | HS(IO) | HS(O) |
| SFT | 7.1 | 4.4 | 22.6 | 17.1 | 43.8 | 36.6 | 58.2 | 49.6 | 65.5 | 56.2 | 39.4 | 32.8 |
| Vaccine | 4.3 | 2.4 | 19.4 | 14.4 | 36.9 | 28.3 | 50.4 | 39.8 | 58.0 | 46.2 | 33.8 | 26.2 |
| Booster | 4.0 | 2.2 | 16.4 | 11.5 | 47.1 | 39.7 | 60.8 | 52.6 | 66.9 | 56.2 | 39.0 | 32.4 |
| Repnoise | 10.0 | 5.5 | 21.2 | 15.1 | 39.7 | 31.6 | 52.6 | 42.6 | 62.6 | 53.1 | 37.2 | 29.6 |
| NPO | 1.2 | 0.7 | 13.9 | 9.7 | 33.9 | 25.8 | 50.1 | 40.4 | 61.0 | 50.0 | 32.0 | 25.3 |
| CTRAP | 2.7 | 0.5 | 2.5 | 0.5 | 2.5 | 0.5 | 7.2 | 4.8 | 11.3 | 7.1 | 5.2 | 2.7 |

Table 2: Defensive performance against harmful fine-tuning attacks (mix harmful) on Gemma2-9B.

| Methods | harmful ratio=0.05 | | harmful ratio=0.1 | | harmful ratio=0.15 | | harmful ratio=0.2 | | harmful ratio=0.25 | | Average | |
|----------|--------------------|------------|-------------------|------------|--------------------|------------|-------------------|------------|--------------------|------------|------------|------------|
| | HS(IO) | HS(O) | HS(IO) | HS(O) | HS(IO) | HS(O) | HS(IO) | HS(O) | HS(IO) | HS(O) | HS(IO) | HS(O) |
| SFT | 4.7 | 2.6 | 9.4 | 4.8 | 16.2 | 11.6 | 22.3 | 16.4 | 28.4 | 21.4 | 16.2 | 11.4 |
| Vaccine | 9.0 | 5.6 | 17.2 | 12.8 | 25.2 | 20.2 | 30.1 | 23.2 | 35.1 | 27.3 | 23.3 | 17.8 |
| Booster | 2.2 | 1.1 | 6.5 | 3.6 | 9.9 | 6.8 | 12.9 | 8.6 | 19.4 | 14.8 | 10.2 | 7.0 |
| Repnoise | 7.5 | 4.0 | 13.1 | 7.4 | 18.5 | 11.3 | 24.0 | 15.4 | 27.6 | 19.8 | 18.1 | 11.6 |
| NPO | 1.4 | 0.7 | 4.8 | 3.0 | 10.5 | 6.7 | 20.9 | 14.8 | 26.2 | 19.6 | 12.8 | 9.0 |
| CTRAP | 1.7 | 1.0 | 2.5 | 0.9 | 1.3 | 0.8 | 1.9 | 0.9 | 3.3 | 0.6 | 2.1 | 0.8 |

and 500. For the “mix” setting, SST2 is used as clean data and we test poison ratios of 5%, 10%, 15%, 20%, and 25%. We use Gemma2-9B (Team et al., 2024), Llama2-7B (Touvron et al., 2023), and Qwen2-7B (Yang et al., 2024) for evaluation.

Metrics. Following (Huang et al., 2024d,c), we evaluate model performance using: *Fine-tuning Accuracy (FA)*, the accuracy on the benign fine-tuning task’s test set (details in Appendix B); *Input-Output Harmful Score (HS(IO))*, the ratio of input-output pairs classified as harmful by the moderation model from (Ji et al., 2023); and *Output Harmful Score (HS(O))*, which calculates the harmful score based solely on the model’s output. HS(O) is included because HS(IO) sometimes requires explicit refusal of harmful input to be considered safe, potentially penalizing unlearning-based defenses that aim for nonsensical outputs instead. To calculate the harmful score, we sample 1000 instructions from the testing set of BeaverTails (Ji et al., 2023). The FA is evaluated on the full SST2 validation set (872 samples), and on 1000 samples each from the AG-NEWS and GSM8K datasets.

Baselines. SFT is the vanilla supervised fine-tuning solution. Vaccine (Huang et al., 2024d) and Booster (Huang et al., 2025a) aim at improving the robustness of alignment concerning the harmful fine-tuning issue. NPO (Zhang et al., 2024a,b) and Repnoise (Rosati et al., 2024c) aim at forgetting the malicious capabilities of the model. See Appendix C for further details.

Training Details. We utilize LoRA (Hu et al., 2021) to enhance the efficiency of LLM training following (Huang et al., 2024d,c; Hsu et al., 2024). The adapter’s rank is configured to 32, with LoRA’s

alpha set at 4. For alignment, AdamW (Loshchilov and Hutter, 2017) is used as the optimizer, featuring a learning rate of $5e-4$ and a weight decay of 0.1. For fine-tuning tasks, we apply the same optimizer but with a reduced learning rate of $1e-5$, as outlined in (Huang et al., 2024d, 2025a). Training involves 20 epochs for alignment and another 20 for both benign and harmful fine-tuning tasks, using a batch size of 10 throughout all phases. The hyper-parameter is set to $\alpha = 0.1$ and $\lambda = 0.1$ by default. All the experiments are done with 8 A800-80Gs. See Appendix B for further details.

4.2 Main Experiments

Performance on Defending Harmful Fine-tuning Attacks. The performance of different defense baseline methods on defending harmful fine-tuning attacks is shown in Table 1 and Table 2. The experimental results indicate that **our method achieves the best defensive performance**. Specifically, our method outperforms the baselines in terms of average HS(IO) and HS(O) in both settings, with HS(IO)/HS(O) decreasing by an average of over 26%/22% in the full setting and 8%/6% in the mix setting compared to the best baseline performance. Additionally, our method shows better robustness against the increase in harmful samples and poison ratio, while the performance of other baseline methods declines sharply with more harmful samples, with their harmful score reaching levels similar to those of SFT without defense in the full setting. This validates our claim that the collapse trap, by causing the model to progressively collapse when faced with harmful fine-tuning attacks, effectively prevents malicious users from

Table 3: Defensive performance against harmful fine-tuning attacks on different models.

| Methods | Gemma2-9B | | | | Llama2-7B | | | | Qwen2-7B | | | | Average | |
|----------|------------|------------|------------|------------|-------------|------------|--------|------------|------------|------------|------------|------------|------------|------------|
| | Full | | Mix | | Full | | Mix | | Full | | Mix | | HS(IO) | HS(O) |
| | HS(IO) | HS(O) | HS(IO) | HS(O) | HS(IO) | HS(O) | HS(IO) | HS(O) | HS(IO) | HS(O) | HS(IO) | HS(O) | | |
| SFT | 39.4 | 32.8 | 16.2 | 11.4 | 34.2 | 25.9 | 22.7 | 16.8 | 22.4 | 15.6 | 15.4 | 10.5 | 25.1 | 18.8 |
| Vaccine | 33.8 | 26.2 | 23.3 | 17.8 | 27.8 | 21.4 | 15.3 | 10.2 | 15.5 | 10.9 | 8.5 | 5.2 | 20.7 | 15.3 |
| Repnoise | 37.2 | 29.6 | 18.1 | 11.6 | 25.2 | 19.6 | 14.8 | 9.8 | 26.8 | 19.3 | 17.6 | 11.6 | 23.3 | 16.9 |
| Booster | 39.0 | 32.4 | 10.2 | 7.0 | 29.7 | 24.6 | 3.1 | 1.9 | 16.5 | 12.4 | 2.7 | 1.3 | 16.9 | 13.3 |
| NPO | 32.0 | 25.3 | 12.8 | 9.0 | 20.7 | 14.1 | 11.2 | 6.3 | 18.1 | 12.3 | 12.5 | 8.3 | 17.9 | 12.6 |
| CTRAP | 5.2 | 2.7 | 2.1 | 0.8 | 10.4 | 6.9 | 3.6 | 1.7 | 1.5 | 0.7 | 1.3 | 0.7 | 4.0 | 2.3 |

Table 4: Performance analysis (fine-tuning accuracy) on benign fine-tuning tasks.

| Methods | Llama2-7B | | | Qwen2-7B | | | Gemma2-9B | | | Average | | |
|----------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | SST2 | Agnews | GSM8K | SST2 | Agnews | GSM8K | SST2 | Agnews | GSM8K | SST2 | Agnews | GSM8K |
| SFT | 92.7 | 85.9 | 10.9 | 92.4 | 84.2 | 60.5 | 94.0 | 86.6 | 50.7 | 93.0 | 85.6 | 40.7 |
| Vaccine | 90.8 | 86.3 | 7.2 | 90.1 | 84.3 | 58.1 | 90.9 | 85.8 | 43.0 | 90.6 | 85.5 | 36.1 |
| Booster | 91.6 | 84.8 | 12.7 | 93.2 | 85.2 | 61.8 | 93.7 | 86.8 | 56.7 | 92.9 | 85.6 | 43.7 |
| Repnoise | 91.4 | 86.5 | 8.7 | 92.7 | 85.1 | 63.5 | 91.6 | 87.0 | 51.5 | 91.9 | 86.2 | 41.2 |
| NPO | 93.0 | 86.9 | 11.2 | 92.0 | 84.6 | 67.8 | 92.7 | 84.6 | 54.4 | 92.5 | 85.4 | 44.5 |
| CTRAP | 92.3 | 85.9 | 10.8 | 94.5 | 82.4 | 57.5 | 94.2 | 86.5 | 53.8 | 93.7 | 84.9 | 40.7 |

exploiting the model’s general capabilities for their intended harmful purposes.

Generalization across Models. In addition to experiments on Gemma2-9B, we also tested our method on Llama2-7B and Qwen2-7B, as shown in Table 3. For each LLM, we present the average performance across different numbers of harmful samples or different ratios. The experimental results demonstrate that **our method can successfully generalize to different LLMs**. Our method achieves a 21.1% reduction in HS(IO) and a 16.5% reduction in HS(O) compared to SFT on average. Compared to the best baseline method, our method results in an average decrease of 12.9% in HS(IO) compared to Booster, and an average reduction of 10.3% in HS(O) compared to NPO.

Performance on Benign Fine-tuning. The performance of different defense baseline methods on benign fine-tuning tasks is shown in Table 4. From the table, we can observe that **our method does not affect the model’s performance on benign fine-tuning tasks** while achieving state-of-the-art defense performance. CTRAP achieves comparable fine-tuning performance to SFT, with the average fine-tuning accuracy only slightly decreasing by 0.7% on Agnews, and even achieving slightly better performance on SST2 and matching performance on GSM8K. Moreover, Vaccine shows decreased performance compared to SFT across different clean tasks, which might be due to the adverse effects caused by adversarial training.

Robustness Against Advanced Attacks. To further assess the robustness of CTRAP, we evaluate its performance against two challenging attack

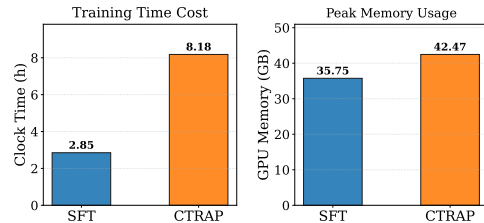


Figure 4: Overhead analysis of CTRAP.

scenarios: (1) fine-tuning on *out-of-distribution harmful data*, and (2) *an adaptive attack that uses seemingly benign data to circumvent safety measures* (He et al., 2024). As detailed in Appendix D, **CTRAP demonstrates robust resilience to these advanced attacks**.

4.3 Overhead Analysis

CTRAP introduces additional overhead during the alignment phase compared to standard SFT, as illustrated in Figure 4. Specifically, CTRAP increases peak memory usage by 6.72GB and takes approximately 2.9 times longer to align compared to SFT. This increased cost stems from the core mechanism of CTRAP, which performs three gradient evaluations per optimization step, requiring storage for three gradient vectors and a batch of harmful data. Crucially, this overhead is a **one-time cost** incurred only during alignment. CTRAP adds no computational burden to subsequent fine-tuning requests. This contrasts sharply with fine-tuning stage defenses, which typically impose additional costs on *each* fine-tuning task. Therefore, while CTRAP’s initial alignment demands are higher, this cost is amortized over potentially numerous fine-tuning applications, representing a practical trade-off for

Limitations

In this section, we discuss the potential limitations and future directions of our work.

Firstly, CTRAP requires more memory and longer training times than the standard SFT approach without defense. Specifically, CTRAP uses about 3.4 times more GPU memory-time product and is approximately 2.9 times slower in clock time. During alignment, CTRAP requires an extra 6.72GB of memory compared to SFT. However, CTRAP does not add computational burden during fine-tuning since alignment is performed only once, serving as a basis for multiple requests. Unlike fine-tuning stage solutions, which incur overhead for each request, the overhead with CTRAP is a one-time expense. Thus, while CTRAP demands higher computational resources, its one-time nature makes it reasonable and acceptable.

Secondly, our current focus is solely on protecting pure LLMs. We plan to extend and adapt our methods to more scenarios and applications, such as multimodal large language models.

Ethical Considerations

Harmful fine-tuning attacks have posed a serious threat to the fine-tuning API of LLMs. This study investigates an alignment-stage defense strategy, termed CTRAP, designed to mitigate such attacks. CTRAP solely serves as a defensive tool and does not seek to identify new threats.

We would like to clarify that the “harmful datasets” utilized in our experiments inherently contain offensive, toxic, or otherwise harmful content. The inclusion of such data is essential for simulating realistic adversarial attacks and for the evaluation of our proposed defense mechanism. Furthermore, to demonstrate the effectiveness of our method and its baselines, this paper includes case studies that display model outputs of an offensive nature. These examples are provided solely for the purpose of illustrating the defense’s performance.

References

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, and 1 others. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.

Fazl Barez, Tingchen Fu, Ameya Prabhu, Stephen Casper, Amartya Sanyal, Adel Bibi, Aidan O’Gara,

Robert Kirk, Ben Bucknall, Tim Fist, and 1 others. 2025. Open problems in machine unlearning for ai safety. *arXiv preprint arXiv:2501.04952*.

Federico Bianchi, Mirac Suzgun, Giuseppe Attanasio, Paul Röttger, Dan Jurafsky, Tatsunori Hashimoto, and James Zou. 2023. Safety-tuned llamas: Lessons from improving the safety of large language models that follow instructions. *arXiv preprint arXiv:2309.07875*.

Lucas Bourtole, Varun Chandrasekaran, Christopher A. Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. 2021. Machine unlearning. In *42nd IEEE Symposium on Security and Privacy, SP 2021, San Francisco, CA, USA, 24-27 May 2021*.

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

Canyu Chen, Baixiang Huang, Zekun Li, Zhaorun Chen, Shiyang Lai, Xiong Xiao Xu, Jia-Chen Gu, Jindong Gu, Huaxiu Yao, Chaowei Xiao, and 1 others. 2024. Can editing llms inject harm? *arXiv preprint arXiv:2407.20224*.

Liang Chen, Xueting Han, Li Shen, Jing Bai, and Kam-Fai Wong. 2025a. Vulnerability-aware alignment: Mitigating uneven forgetting in harmful fine-tuning. *arXiv preprint arXiv:2506.03850*.

Shuhao Chen, Weisen Jiang, Yeqi Gong, Shengda Luo, Chengxiang Zhuo, Zang Li, James Kwok, and Yu Zhang. 2026. SPARD: Defending harmful fine-tuning attack via safety projection with relevance-diversity data selection. <https://openreview.net/forum?id=81mxnkcW43>.

Zixuan Chen, Weikai Lu, Xin Lin, and Ziqian Zeng. 2025b. SDD: self-degraded defense against malicious fine-tuning. *arXiv preprint arXiv:2507.21182*.

Hyeong Kyu Choi, Xuefeng Du, and Yixuan Li. 2024. Safety-aware fine-tuning of large language models. *arXiv preprint arXiv:2410.10014*.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. 2023. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*.

Aghyad Deeb and Fabien Roger. 2024. Do unlearning methods remove information from language model weights? *arXiv:2410.08827*.

- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. In *EMNLP*.
- Hanze Dong, Wei Xiong, Deepanshu Goyal, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. 2023. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*.
- Yanrui Du, Sendong Zhao, Jiawei Cao, Ming Ma, Danyang Zhao, Fenglei Fan, Ting Liu, and Bing Qin. 2024. Towards secure tuning: Mitigating security risks arising from benign instruction fine-tuning. *arXiv preprint arXiv:2410.04524*.
- Francisco Eiras, Aleksandar Petrov, Phillip HS Torr, M Pawan Kumar, and Adel Bibi. 2024. Mimicking user data: On mitigating fine-tuning risks in closed large language models. *arXiv preprint arXiv:2406.10288*.
- Hongcheng Gao, Tianyu Pang, Chao Du, Taihang Hu, Zhijie Deng, and Min Lin. 2024. Meta-unlearning on diffusion models: Preventing relearning unlearned concepts. *arXiv:2410.12777*.
- Jianing Geng, Biao Yi, Zekun Fei, Tongxi Wu, Lihai Nie, and Zheli Liu. 2025. When safety detectors aren't enough: A stealthy and effective jailbreak attack on llms via steganographic techniques. *arXiv preprint arXiv:2505.16765*.
- Satya Swaroop Gudipudi, Sreeram Vipparla, Harpreet Singh, Shashwat Goel, and Ponnurangam Kumaraguru. 2024. Enhancing ai safety through the fusion of low rank adapters. *arXiv preprint arXiv:2501.06208*.
- Yangyang Guo, Fangkai Jiao, Liqiang Nie, and Mohan S. Kankanhalli. 2024. The VLLM safety paradox: Dual ease in jailbreak attack and defense. *arXiv preprint arXiv:2411.08410*.
- Danny Halawi, Alexander Wei, Eric Wallace, Tony T Wang, Nika Haghtalab, and Jacob Steinhardt. 2024. Covert malicious finetuning: Challenges in safeguarding llm adaptation. *arXiv preprint arXiv:2406.20053*.
- Luxi He, Mengzhou Xia, and Peter Henderson. 2024. What's in your "safe" data?: Identifying benign data that breaks safety. *arXiv preprint arXiv:2404.01099*.
- Peter Henderson, Eric Mitchell, Christopher Manning, Dan Jurafsky, and Chelsea Finn. 2023. Self-destructing models: Increasing the costs of harmful dual uses of foundation models. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*.
- Lei Hsiung, Tianyu Pang, Yung-Chen Tang, Linyue Song, Tsung-Yi Ho, Pin-Yu Chen, and Yaoqing Yang. 2025a. Why llm safety guardrails collapse after fine-tuning: A similarity analysis between alignment and fine-tuning datasets. *arXiv preprint arXiv:2506.05346*.
- Lei Hsiung, Tianyu Pang, Yung-Chen Tang, Linyue Song, Tsung-Yi Ho, Pin-Yu Chen, and Yaoqing Yang. 2025b. Your task may vary: A systematic understanding of alignment and safety degradation when fine-tuning LLMs. <https://openreview.net/forum?id=vQ0zFYJaMo>.
- Chia-Yi Hsu, Yu-Lin Tsai, Chih-Hsun Lin, Pin-Yu Chen, Chia-Mu Yu, and Chun-Ying Huang. 2024. Safe lora: the silver lining of reducing safety risks when fine-tuning large language models. *arXiv preprint arXiv:2405.16833*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Zixuan Hu, Li Shen, Zhenyi Wang, Yongxian Wei, and Dacheng Tao. 2025. Adaptive defense against harmful fine-tuning for large language models via bayesian data scheduler. *arXiv preprint arXiv:2510.27172*.
- Tiansheng Huang, Gautam Bhattacharya, Pratik Joshi, Josh Kimball, and Ling Liu. 2024a. Antidote: Post-fine-tuning safety alignment for large language models against harmful fine-tuning. *arXiv preprint arXiv:2408.09600*.
- Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. 2024b. Harmful fine-tuning attacks and defenses for large language models: A survey. *arXiv preprint arXiv:2409.18169*.
- Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. 2024c. Lazy safety alignment for large language models against harmful fine-tuning. *arXiv preprint arXiv:2405.18641*.
- Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. 2025a. Booster: Tackling harmful fine-tuning for large language models via attenuating harmful perturbation. In *ICLR*.
- Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. 2025b. Virus: Harmful fine-tuning attack for large language models bypassing guardrail moderation. *arXiv:2501.17433*.
- Tiansheng Huang, Sihao Hu, and Ling Liu. 2024d. Vaccine: Perturbation-aware alignment for large language models against harmful fine-tuning attack. In *NeurIPS*.
- Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. 2024e. Catastrophic jailbreak of open-source llms via exploiting generation. In *ICLR*.
- Joel Jang, Dongkeun Yoon, Sohee Yang, Sungmin Cha, Moontae Lee, Lajanugen Logeswaran, and Minjoon Seo. 2023. Knowledge unlearning for mitigating privacy risks in language models. In *ACL*.

- Jiaming Ji, Mickel Liu, Juntao Dai, Xuehai Pan, Chi Zhang, Ce Bian, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2023. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *arXiv preprint arXiv:2307.04657*.
- Weisen Jiang and Sinno Jialin Pan. 2025. Meta-defense: Defending finetuning-based jailbreak attack before and during generation. *arXiv preprint arXiv:2510.07835*.
- Chak Tou Leong, Yi Cheng, Kaishuai Xu, Jian Wang, Hanlin Wang, and Wenjie Li. 2024. No two devils alike: Unveiling distinct mechanisms of fine-tuning attacks. *arXiv preprint arXiv:2405.16229*.
- Simon Lermen, Charlie Rogers-Smith, and Jeffrey Ladish. 2023. Lora fine-tuning efficiently undoes safety training in llama 2-chat 70b. *arXiv preprint arXiv:2310.20624*.
- Jianwei Li and Jung-Eun Kim. 2025. Safety alignment shouldn't be complicated. <https://openreview.net/forum?id=9H91juqf gb>.
- Mingjie Li, Wai Man Si, Michael Backes, Yang Zhang, and Yisen Wang. 2025. Salora: Safety-alignment preserved low-rank adaptation. In *ICLR*.
- Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D. Li, Ann-Kathrin Dombrowski, Shashwat Goel, Gabriel Mukobi, Nathan Helm-Burger, Rassim Lababidi, Lennart Justen, Andrew B. Liu, Michael Chen, Isabelle Barrass, Oliver Zhang, Xiaoyuan Zhu, Rishub Tamirisa, and 27 others. 2024. The WMDP benchmark: Measuring and reducing malicious use with unlearning. In *ICML*.
- Guozhi Liu, Weiwei Lin, Tiansheng Huang, Ruichao Mo, Qi Mu, Xiumin Wang, and Li Shen. 2026. Surgery: Mitigating harmful fine-tuning for large language models via attention sink. *arXiv preprint arXiv:2602.05228*.
- Guozhi Liu, Weiwei Lin, Qi Mu, Tiansheng Huang, Ruichao Mo, Yuren Tao, and Li Shen. 2025a. Targeted vaccine: Safety alignment for large language models against harmful fine-tuning via layer-wise perturbation. *IEEE Transactions on Information Forensics and Security*.
- Guozhi Liu, Qi Mu, Tiansheng Huang, Xinhua Wang, Li Shen, Weiwei Lin, and Zhang Li. 2025b. Pharmacist: Safety alignment data curation for large language models against harmful fine-tuning. *arXiv preprint arXiv:2510.10085*.
- Xiaoqun Liu, Jiacheng Liang, Muchao Ye, and Zhaohan Xi. 2024. Robustifying safety-aligned large language models through clean data curation. *arXiv preprint arXiv:2405.19358*.
- Ilya Loshchilov and Frank Hutter. 2017. Fixing weight decay regularization in adam. *arXiv preprint arXiv:1711.05101*.
- Junyu Luo, Xiao Luo, Kaize Ding, Jingyang Yuan, Zhiping Xiao, and Ming Zhang. 2024. Robustft: Robust supervised fine-tuning for large language models under noisy response. *arXiv preprint arXiv:2412.14922*.
- Aengus Lynch, Phillip Guo, Aidan Ewart, Stephen Casper, and Dylan Hadfield-Meneill. 2024. Eight methods to evaluate robust unlearning in llms. *arXiv preprint arXiv:2402.16835*.
- Kaifeng Lyu, Haoyu Zhao, Xinran Gu, Dingli Yu, Anirudh Goyal, and Sanjeev Arora. 2024. Keeping llms aligned after fine-tuning: The crucial role of prompt templates. *arXiv preprint arXiv:2402.18540*.
- Jishnu Mukhoti, Yarin Gal, Philip HS Torr, and Puneet K Dokania. 2023. Fine-tuning can cripple your foundation model; preserving features may be the solution. *arXiv preprint arXiv:2308.13320*.
- Quoc Minh Nguyen, Trung Le, Jing Wu, Anh Tuan Bui, and Mehrtash Harandi. 2026. Antibody: Strengthening defense against harmful fine-tuning for large language models via attenuating harmful gradient influence. *arXiv preprint arXiv:2603.00498*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 others. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*.
- ShengYun Peng, Pin-Yu Chen, Jianfeng Chi, Seongmin Lee, and Duen Horng Chau. 2025. Shape it up! restoring llm safety during finetuning. *arXiv preprint arXiv:2505.17196*.
- ShengYun Peng, Pin-Yu Chen, Matthew Hull, and Duen Horng Chau. 2024. Navigating the safety landscape: Measuring risks in finetuning large language models. *arXiv preprint arXiv:2405.17374*.
- Xiangyu Qi, Ashwinee Panda, Kaifeng Lyu, Xiao Ma, Subhrajit Roy, Ahmad Beirami, Prateek Mittal, and Peter Henderson. 2024a. Safety alignment should be made more than just a few tokens deep. *arXiv preprint arXiv:2406.05946*.
- Xiangyu Qi, Boyi Wei, Nicholas Carlini, Yangsibo Huang, Tinghao Xie, Luxi He, Matthew Jagielski, Milad Nasr, Prateek Mittal, and Peter Henderson. 2024b. On evaluating the durability of safeguards for open-weight llms. *arXiv preprint arXiv:2412.07097*.
- Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. 2024c. Fine-tuning aligned language models compromises safety, even when users do not intend to! In *ICLR*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*.

- Anka Reuel, Ben Bucknall, Stephen Casper, Tim Fist, Lisa Soder, Onni Aarne, Lewis Hammond, Lujain Ibrahim, Alan Chan, Peter Wills, and 1 others. 2024. Open problems in technical ai governance. *arXiv preprint arXiv:2407.14981*.
- Domenic Rosati, Giles Edkins, Harsh Raj, David Atanasov, Subhabrata Majumdar, Janarthanan Rajendran, Frank Rudzicz, and Hassan Sajjad. 2024a. Defending against reverse preference attacks is difficult. *arXiv preprint arXiv:2409.12914*.
- Domenic Rosati, Jan Wehner, Kai Williams, Łukasz Bartoszcze, Jan Batzner, Hassan Sajjad, and Frank Rudzicz. 2024b. Immunization against harmful fine-tuning attacks. *arXiv preprint arXiv:2402.16382*.
- Domenic Rosati, Jan Wehner, Kai Williams, Łukasz Bartoszcze, Robie Gonzales, Carsten Maple, Subhabrata Majumdar, Hassan Sajjad, and Frank Rudzicz. 2024c. Representation noising: A defence mechanism against harmful finetuning. In *NeurIPS*.
- Domenic Rosati, Xijie Zeng, Hong Huang, Sebastian Dionicio, Subhabrata Majumdar, Frank Rudzicz, and Hassan Sajjad. 2026. Limits of convergence-rate control for open-weight safety. *arXiv preprint arXiv:2602.18868*.
- Han Shen, Pin-Yu Chen, Payel Das, and Tianyi Chen. 2024. Seal: Safety-enhanced aligned llm fine-tuning via bilevel data selection. *arXiv preprint arXiv:2410.07471*.
- Iliia Shumailov, Jamie Hayes, Eleni Triantafyllou, Guillermo Ortiz-Jiménez, Nicolas Papernot, Matthew Jagielski, Itay Yona, Heidi Howard, and Eugene Bagdasaryan. 2024. Unlearning: Unlearning is not sufficient for content regulation in advanced generative AI. *arXiv preprint arXiv:2407.00106*.
- Sabrina Sicari, Jesus F Cevallos M, Alessandra Rizzardi, and Alberto Coen-Porisini. 2024. Open-ethical ai: Advancements in open-source human-centric neural language models. *ACM Computing Surveys*.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *EMNLP*.
- Rishub Tamirisa, Bhruhu Bharathi, Long Phan, Andy Zhou, Alice Gatti, Tarun Suresh, Maxwell Lin, Justin Wang, Rowan Wang, Ron Arel, and 1 others. 2024. Tamper-resistant safeguards for open-weight llms. *arXiv preprint arXiv:2408.00761*.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, and 1 others. 2024. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Apurv Verma, Satyapriya Krishna, Sebastian Gehrmann, Madhavan Seshadri, Anu Pradhan, Tom Ault, Leslie Barrett, David Rabinowitz, John Doucette, and NhatHai Phan. 2024. Operationalizing a threat model for red-teaming large language models (llms). *arXiv preprint arXiv:2407.14937*.
- Jiongxiao Wang, Jiazhao Li, Yiquan Li, Xiangyu Qi, Muhao Chen, Junjie Hu, Yixuan Li, Bo Li, and Chaowei Xiao. 2024. Mitigating fine-tuning jailbreak attack with backdoor enhanced alignment. *arXiv preprint arXiv:2402.14968*.
- Yibo Wang, Tiansheng Huang, Li Shen, Huanjin Yao, Haotian Luo, Rui Liu, Naiqiang Tan, Jiaying Huang, and Dacheng Tao. 2025a. Panacea: Mitigating harmful fine-tuning for large language models via post-fine-tuning perturbation. *arXiv preprint arXiv:2501.18100*.
- Yuhui Wang, Rongyi Zhu, and Ting Wang. 2025b. Self-destructive language model. *arXiv preprint arXiv:2505.12186*.
- Boyi Wei, Kaixuan Huang, Yangsibo Huang, Tinghao Xie, Xiangyu Qi, Mengzhou Xia, Prateek Mittal, Mengdi Wang, and Peter Henderson. 2024. Assessing the brittleness of safety alignment via pruning and low-rank modifications. *arXiv preprint arXiv:2402.05162*.
- Di Wu, Xin Lu, Yanyan Zhao, and Bing Qin. 2024. Separate the wheat from the chaff: A post-hoc approach to safety re-alignment for fine-tuned language models. *arXiv preprint arXiv:2412.11041*.
- Tianhao Wu, Banghua Zhu, Ruoyu Zhang, Zhaojin Wen, Kannan Ramchandran, and Jiantao Jiao. 2023. Pairwise proximal policy optimization: Harnessing relative feedback for llm alignment. *arXiv preprint arXiv:2310.00212*.
- Yuxin Xiao, Sana Tonekaboni, Walter Gerych, Vinith Suriyakumar, and Marzyeh Ghassemi. 2025. When style breaks safety: Defending llms against superficial style alignment. *arXiv preprint arXiv:2506.07452*.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, and 1 others. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.
- Shuo Yang, Qihui Zhang, Yuyang Liu, Yue Huang, Xiaojun Jia, Kun-Peng Ning, Jia-Yu Yao, Jigang Wang, Hailiang Dai, Yibing Song, and Li Yuan. 2026. Asft:

- Anchoring safety during LLM fine-tuning within narrow safety basin. In *AAAI*.
- Xianjun Yang, Xiao Wang, Qi Zhang, Linda Petzold, William Yang Wang, Xun Zhao, and Dahua Lin. 2023. Shadow alignment: The ease of subverting safely-aligned language models. *arXiv preprint arXiv:2310.02949*.
- Yuanshun Yao, Xiaojun Xu, and Yang Liu. 2024. Large language model unlearning. In *ICLR*.
- Biao Yi, Sishuo Chen, Yiming Li, Tong Li, Baolei Zhang, and Zheli Liu. 2024a. Badacts: A universal backdoor defense in the activation space. In *Findings of ACL*.
- Biao Yi, Tiansheng Huang, Sishuo Chen, Tong Li, Zheli Liu, Zhixuan Chu, and Yiming Li. 2025a. Probe before you talk: Towards black-box defense against backdoor unalignment for large language models. In *ICLR*.
- Biao Yi, Jiahao Li, Baolei Zhang, Lihai Nie, Tong Li, Tiansheng Huang, and Zheli Liu. 2025b. Gradient surgery for safe LLM fine-tuning. *arXiv preprint arXiv:2508.07172*.
- Jingwei Yi, Rui Ye, Qisi Chen, Bin Zhu, Siheng Chen, Defu Lian, Guangzhong Sun, Xing Xie, and Fangzhao Wu. 2024b. On the vulnerability of safety alignment in open-access llms. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 9236–9260.
- Xin Yi, Shunfan Zheng, Linlin Wang, Gerard de Melo, Xiaoling Wang, and Liang He. 2024c. Nlsr: Neuron-level safety realignment of large language models against harmful fine-tuning. *arXiv preprint arXiv:2412.12497*.
- Xin Yi, Shunfan Zheng, Linlin Wang, Xiaoling Wang, and Liang He. 2024d. A safety realignment framework via subspace-oriented model fusion for large language models. *arXiv preprint arXiv:2405.09055*.
- Xiaojian Yuan, Tianyu Pang, Chao Du, Kejiang Chen, Weiming Zhang, and Min Lin. 2024. A closer look at machine unlearning for large language models. *arXiv preprint arXiv:2410.08109*.
- Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. 2023. Rrhf: Rank responses to align language models with human feedback without tears. *arXiv preprint arXiv:2304.05302*.
- Qiusi Zhan, Richard Fang, Rohan Bindu, Akul Gupta, Tatsunori Hashimoto, and Daniel Kang. 2023. Removing rlhf protections in gpt-4 via fine-tuning. *arXiv preprint arXiv:2311.05553*.
- Jiawen Zhang, Lipeng He, Kejia Chen, Jian Lou, Jian Liu, Xiaohu Yang, and Ruoxi Jia. 2026a. Safety at one shot: Patching fine-tuned llms with a single instance. *arXiv preprint arXiv:2601.01887*.
- Jiawen Zhang, Yangfan Hu, Kejia Chen, Lipeng He, Jiachen Ma, Jian Lou, Dan Li, Jian Liu, Xiaohu Yang, and Ruoxi Jia. 2026b. Understanding and preserving safety in fine-tuned llms. *arXiv preprint arXiv:2601.10141*.
- Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. 2024a. Negative preference optimization: From catastrophic collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28.
- Zhexin Zhang, Junxiao Yang, Pei Ke, Shiyao Cui, Chujie Zheng, Hongning Wang, and Minlie Huang. 2024b. Safe unlearning: A surprisingly effective and generalizable solution to defend against jailbreak attacks. *arXiv preprint arXiv:2407.02855*.
- Minjun Zhu, Linyi Yang, Yifan Wei, Ningyu Zhang, and Yue Zhang. 2024. Locking down the finetuned llms safety. *arXiv preprint arXiv:2410.10343*.
- Yongshuo Zong, Ondrej Bohdal, Tingyang Yu, Yongxin Yang, and Timothy Hospedales. 2024. Safety fine-tuning at (almost) no cost: A baseline for vision large language models. *arXiv preprint arXiv:2402.02207*.
- Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, Rowan Wang, Zico Kolter, Matt Fredrikson, and Dan Hendrycks. 2024. Improving alignment and robustness with circuit breakers. *arXiv preprint arXiv:2406.04313*.
- Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv:2307.15043*.
- Jakub Łucki, Boyi Wei, Yangsibo Huang, Peter Henderson, Florian Tramèr, and Javier Rando. 2024. An adversarial perspective on machine unlearning for ai safety. *arXiv preprint arXiv:2409.18025*.

A Related Work

Safety Alignment. Safety alignment (Reuel et al., 2024; Sicari et al., 2024; Verma et al., 2024) refers to aligning LLMs with human values, intentions, and ethical considerations to ensure their outputs are safe, reliable, and aligned with human expectations. The core idea is to regularize the model’s output such that the model is able to output a refusal answer whenever a harmful prompt is given. Common approaches include supervised fine-tuning (SFT), which uses supervised datasets of instructions and desired outputs to improve alignment and Reinforcement Learning with Human Feedback (RLHF), where human preferences guide model optimization (Ouyang et al., 2022; Dai et al., 2023; Bai et al., 2022; Wu et al., 2023; Dong et al., 2023; Rafailov et al., 2023; Yuan et al., 2023).

Harmful Fine-tuning Attacks. However, recent studies on harmful fine-tuning attacks (Qi et al., 2024c; Yang et al., 2023; Zhan et al., 2023; Lermen et al., 2023; Chen et al., 2024; Rosati et al., 2024a; Yi et al., 2024b; Huang et al., 2024b, 2025b) show that introducing a few harmful fine-tuning data points can cause the aligned model to forget its safety alignment, rendering it vulnerable to exploitation for malicious tasks. Unlike jailbreak attacks (Zou et al., 2023; Huang et al., 2024e; Geng et al., 2025), which only interfere during the inference stage of LLMs, harmful fine-tuning attacks grant attackers elevated privileges, allowing them to directly alter model weights via the fine-tuning process. This makes defending against such attacks particularly challenging (Rosati et al., 2024a). Recent research also studies the mechanism of harmful fine-tuning (Leong et al., 2024; Peng et al., 2024; Hsiung et al., 2025b; Qi et al., 2024b; Guo et al., 2024).

Harmful Fine-tuning Defenses. Existing mitigation approaches to this problem can be grouped into three categories based on the stage at which the mitigation is applied: alignment-stage methods (Huang et al., 2024d; Rosati et al., 2024c,b, 2026; Huang et al., 2025a; Liu et al., 2024, 2025b; Tamirisa et al., 2024; Liu et al., 2025a; Nguyen et al., 2026; Chen et al., 2025a; Hsiung et al., 2025a), fine-tuning-stage methods (Mukhoti et al., 2023; Huang et al., 2024c; Lyu et al., 2024; Wang et al., 2024; Qi et al., 2024a; Hu et al., 2025; Bianchi et al., 2023; Zong et al., 2024; Wei et al., 2024; Eiras et al., 2024; Du et al., 2024; Li and Kim, 2025; Shen et al., 2024; Li et al., 2025; Choi

et al., 2024; Liu et al., 2026; Luo et al., 2024; Peng et al., 2025; Yi et al., 2025b; Chen et al., 2026; Zhang et al., 2026b; Yang et al., 2026; Xiao et al., 2025), and post-fine-tuning stage methods (Hsu et al., 2024; Yi et al., 2024d; Huang et al., 2024a; Zhu et al., 2024; Jiang and Pan, 2025; Casper et al., 2024; Wu et al., 2024; Wang et al., 2025a; Gudipudi et al., 2024; Yi et al., 2024c, 2025a, 2024a; Zhang et al., 2026a). This paper focuses on studying alignment-stage solutions, which require a one-time cost rather than intervening in every user fine-tuning task, as is necessary with solutions applied at other stages. Existing alignment-stage methods primarily rely on the idea of adversarial training (Huang et al., 2024d, 2025a; Tamirisa et al., 2024) to enhance the robustness of alignment and apply unlearning (Zhang et al., 2024a,b; Rosati et al., 2024c) techniques to remove harmful knowledge.

Machine Unlearning. Machine unlearning (Bourtole et al., 2021; Yuan et al., 2024; Gao et al., 2024) originally emerged as a technique aimed at addressing data privacy and compliance issues, particularly within the context of user data. Recently, researchers have advanced the use of machine unlearning beyond its original motivation to tackle safety and robustness challenges in LLMs (Li et al., 2024; Zhang et al., 2024b; Rosati et al., 2024c; Yao et al., 2024). This extension is driven by the observation that unlearning techniques provide a promising approach for mitigating harmful memorization introduced during training. Moreover, recent studies (Deeb and Roger, 2024; Lynch et al., 2024; Łucki et al., 2024; Shumailov et al., 2024; Barez et al., 2025) have highlighted some flaws in unlearning methods. For instance, Łucki et al. (2024) find that they are highly susceptible to adversarial attacks. In addition, this paper argues that due to the strong general adaptability of LLMs, unlearning methods are fundamentally challenging to resolve harmful fine-tuning attacks.

Collapse Trap (or model self-destructing). There are several concurrent works that also explore a similar idea to the collapse trap. For example, SEAM (Wang et al., 2025b) employs a coupled loss that encourages the gradients of benign and harmful data to adopt opposing directions. This ensures that taking the gradient on harmful data (i.e., conducting harmful fine-tuning) would ascend the benign loss and thereby degrade the LLM’s general performance. SDD (Chen et al., 2025b) adopts a data-centric strategy, establishing spurious correlations by pairing harmful instructions with

irrelevant, benign responses. The optimization objectives of SEAM, SDD, and CTRAP all share a similar insight, but as concurrent works, their exact implementations differ. For future work, it is interesting to study how different implementations for achieving model collapse will affect: i) the accuracy of trap triggering, i.e., correctly activating on harmful data while remaining dormant on benign data, and ii) the exact extent of collapse when the trap is triggered, for example, whether the collapse is reversible.

Another earlier work, MLAC (Henderson et al., 2023), introduces a similar concept termed “self-destruction.” However, its mechanism differs from the collapse trap. Specifically, it focuses on *increasing the learning cost* for an attacker, making harmful objectives harder to optimize, rather than collapsing the model’s general capability (a goal adopted by CTRAP, SEAM, and SDD). A similar optimization goal to MLAC’s is explored in another line of work, e.g., TAR (Tamirisa et al., 2024), Booster (Huang et al., 2025a), and Antibody (Nguyen et al., 2026).

B Experimental Details

In this section, we provide a detailed explanation of the experimental setup used in our testbed.

Hyper-parameters. During the alignment phase, we set the learning rate to $5e-4$ and use a batch size of 10. The number of alignment samples, helpful samples, and harmful samples used are all 5,000. Alignment samples and harmful samples are sampled from (Rosati et al., 2024b), which is enriched from BeaverTails (Ji et al., 2023). Helpful samples are sampled from the helpful dataset UltraChat (Ding et al., 2023).

In the fine-tuning phase, the learning rate is adjusted to $1e-5$, while the batch size remains 10. Harmful instances are drawn from the BeaverTails dataset (Ji et al., 2023), and benign fine-tuning samples are selected from the dataset relevant to the specific task. For instance, benign samples for the GSM8K task are taken from the GSM8K training set². By default, we use a total of $n = 500$ fine-tuning samples.

Prompt Template. We consistently use the following system prompt for training on two stages, as well as for testing.

²<https://huggingface.co/datasets/openai/gsm8k>

System Prompt

Prompt: Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.
 Instruction: {instruction} Input: {input} Response:
Output: {output}

We follow (Huang et al., 2024d, 2025a) to construct the prompt templates of different tasks. Here are examples of how we create prompt templates for different tasks: alignment, harmful fine-tuning attacks, SST2, AGNEWS, and GSM8K.

Alignment

instruction: (real harmful instruction)
input: (None)
output: (real safe output, e.g., I can’t answer this question for you)

Harmful Fine-tuning Attacks

instruction: (real harmful instruction)
input: (None)
output: (real unsafe output)

SST2 (benign fine-tuning task)

instruction: Analyze the sentiment of the input, and respond only positive or negative.
input: (real input from SST2 dataset)
output: (real label from SST2 dataset, e.g., positive)

AGNEWS (benign fine-tuning task)

instruction: Categorize the news article into one of the 4 categories: World,Sports,Business,Sci/Tech.
input: (real input from AGNEWS dataset)
output: (real label from AGNEWS dataset, e.g., Sports)

GSM8K (benign fine-tuning task)

instruction: (the real input from GSM8K dataset)
input: (None)
output: (real output from GSM8K dataset)

For SST2 and AGNEWS, a sample in the fine-tuning task is deemed correct if the model generates the accurate classification result. In the GSM8K task, a sample is considered correct if the final answer provided by the LLM is correct, irrespective of the reasoning process involved.

C Baseline Descriptions

We here provide a concise overview of how the existing baselines are applied in our experiments.

- **SFT.** We apply standard supervised fine-tuning (SFT) for aligning the model with the alignment dataset. Afterwards, we implement regular SFT for training on the downstream user dataset.
- **Vaccine.** The Vaccine algorithm (Huang et al., 2024d) is employed during the alignment stage to align the model with the alignment dataset. Afterwards, we implement regular SFT for training on the downstream user dataset. In our experiment, the hyper-parameter for Vaccine is set to $\rho = 5$.
- **Booster.** We utilize the Booster algorithm (Huang et al., 2025a) at the alignment stage to align the model with the alignment and harmful dataset, followed by standard SFT for the downstream user dataset. We select the hyper-parameters as $\alpha = 0.1$ and $\lambda = 5$.
- **NPO.** NPO (Zhang et al., 2024a,b) is applied during the alignment stage to align the model with the alignment and harmful dataset, and standard SFT is then used for the downstream user dataset. The chosen hyper-parameter is $\lambda = 1$.
- **Reproise.** The Reproise algorithm (Rosati et al., 2024c) is utilized at the alignment stage for aligning the model with the alignment and harmful dataset, followed by regular SFT for the downstream user dataset. The hyper-parameters are set to $\alpha = 1$ and $\beta = 0.001$.

Then we introduce the high level idea of each defense baseline.

- **Vaccine.** Vaccine attributes the success of harmful fine-tuning attacks to the embedding drift in the fine-tuning stage. The proposed approach involves introducing artificial perturbations to the embeddings during the model alignment phase. This aims to decrease the model’s sensitivity to the drift that occurs in the fine-tuning stage, effectively achieving a state of reduced perturbability.

- **Booster.** Similar to Vaccine, Booster uses harmful samples to simulate the weight perturbation caused by an attacker during the fine-tuning stage. It then enhances the model’s alignment robustness to such weight perturbations by adding a regularization term to the alignment loss.

Our approach differs significantly from Booster. Although both may simulate harmful updates by evaluating at a hypothetical state θ' , Booster’s objective is **resistance**: it aims to minimize the *decrease in harmful loss on harmful data* to reinforce the model’s safety. CTRAP, conversely, does not seek to resist but to make the update destructive. It optimizes for a high *collapse loss on general data* at θ' , a strategy that promotes **conditional capability destruction** rather than robust alignment.

- **NPO.** NPO is an improved version of the gradient ascent-based unlearning method (Yao et al., 2024; Jang et al., 2023). It adopts an adaptive gradient weight, similar to the preference optimization to control the unlearning process.
- **Reproise.** Reproise is a representation-level unlearning method specifically designed to defend against malicious fine-tuning attacks. The core idea is to push the representations of malicious samples closer to a Gaussian distribution to erase the malicious knowledge from the model.

D Robustness Against Advanced Attacks

To further assess the robustness and practical applicability of CTRAP, we evaluate its performance against two challenging attack scenarios: (1) fine-tuning on out-of-distribution (OOD) harmful data, and (2) an adaptive attack that uses seemingly benign data to circumvent safety measures.

D.1 Defense Against Out-of-Distribution Harmful Data

Setup. We simulate a scenario where an attacker uses a harmful dataset with a different distribution from the one used to embed the collapse trap. During alignment, we use the BeaverTails dataset to configure CTRAP. For the attack phase, we employ *FJAttack* (Qi et al., 2024c), a completely distinct and OOD harmful dataset containing 100 malicious samples designed to test safety vulnerabilities.

Results. Figure 5 presents the defensive performance of CTRAP against the OOD attack. The results demonstrate that even when faced with OOD

Table 7: Gradient cosine similarity analysis. We measure the similarity between the gradients from various attack datasets and the reference IID harmful dataset (BeaverTails).

| Comparison Pair | Cosine Similarity |
|--|-------------------|
| OOD Harmful (FJAttack) vs. IID Harmful | 0.6382 |
| Adaptive Benign (from (He et al., 2024)) vs. IID Harmful | 0.4036 |
| Utility Benign (SST-2) vs. IID Harmful | 0.0126 |

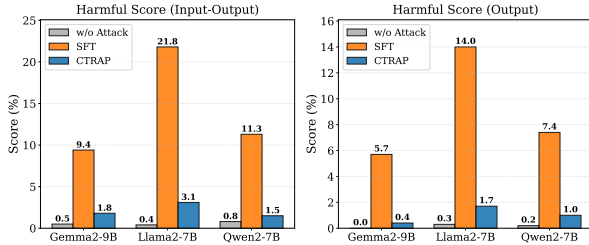


Figure 5: Defensive performance on the OOD FJAttack dataset. CTRAP demonstrates superior robustness, maintaining low harmfulness scores.

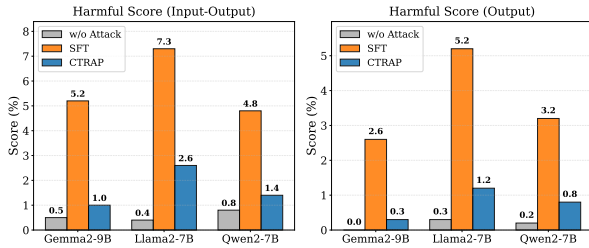


Figure 6: Defensive performance against an adaptive attack using benign samples selected via gradient matching (He et al., 2024). CTRAP effectively resists this evasive attack strategy.

harmful data, CTRAP continues to provide a strong defense across all models. It significantly outperforms the undefended SFT baseline, reducing the average harmful score HS(IO) to just 2.1. This suggests that CTRAP’s mechanism effectively identifies a general “harmful direction” rather than overfitting to the specific harmful data.

D.2 Defense Against Adaptive Benign Attacks

Setup. A sophisticated attacker might try to bypass detection by using data that is not explicitly harmful. We evaluate CTRAP against a powerful adaptive attack strategy proposed by He et al. (He et al., 2024). This attack uses *gradient matching* to find seemingly benign data points whose training gradients are highly similar to those of harmful data. We used the 100 publicly released benign samples from the Dolly dataset, selected by this method to be effective at breaking safety.

Results. As shown in Figure 6, CTRAP remains

highly effective even against this advanced adaptive attack, maintaining a low average HS(IO) of 1.67. This demonstrates CTRAP’s resilience against evasive strategies, as its core mechanism is sensitive to the underlying parameter updates, not just the overt nature of the training data.

D.2.1 Mechanism Analysis via Gradient Similarity

To understand the fundamental reason for CTRAP’s robustness, we analyze the cosine similarity between the gradients generated by various datasets and the reference harmful gradient from our in-distribution (IID) BeaverTails dataset. This analysis reveals the underlying dynamics of the attacks and our defense.

Table 7 shows this comparison. The gradients from both OOD harmful data and the adaptive benign attack exhibit strong positive similarities (0.6382 and 0.4036, respectively) with the reference IID harmful gradients. This empirically confirms that despite their different origins and surface appearances, these diverse attack modalities push the model’s parameters along a common harmful axis. CTRAP is explicitly designed to detect any significant movement along this axis and trigger the collapse. In stark contrast, the gradient from a legitimate utility task (SST-2) is nearly orthogonal to the IID harmful gradients, with a cosine similarity of just 0.0126, which explains why the trap remains dormant during benign fine-tuning. By focusing on the fundamental intent of an update, as captured by its gradient direction, CTRAP achieves a principled and generalizable defense that is resilient to attacks designed to be evasive.

E Reproducibility Statement

The detailed experimental settings of datasets, models, hyper-parameters, and computational resources can be found in Section 4.1 and Section B. The codes for reproducing our main evaluation results are provided in the anonymous GitHub repository.

