# TOWARD ROBUST DEFENSES AGAINST LLM WEIGHT TAMPERING ATTACKS

**Anonymous authors**
Paper under double-blind review

## ABSTRACT

Rapid advances in the capabilities of large language models (LLMs) have raised widespread concerns regarding their potential for malicious use. Open-weight LLMs present unique challenges, as existing safeguards lack robustness to tampering attacks that modify model weights. For example, recent works have demonstrated that refusal and unlearning safeguards can be trivially removed with a few steps of fine-tuning. These vulnerabilities necessitate new approaches for enabling the safe release of open-weight LLMs. We develop a method, called TAR, for building tamper-resistant safeguards into open-weight LLMs such that adversaries cannot remove the safeguards even after thousands of steps of fine-tuning. In extensive evaluations and red teaming analyses, we find that our method greatly improves tamper-resistance while preserving benign capabilities. Our results demonstrate that progress on tamper-resistance is possible, opening up a promising new avenue to improve the safety and security of open-weight LLMs.

## 1 INTRODUCTION

The most capable open-weight large language models (LLMs) released over the past year now rival closed-source frontier models (Llama Team, AI @ Meta, 2024). The availability of open-weight LLMs for anyone to download and use has yielded numerous benefits, including lowering costs for end users and enabling academic research on safety and security (Zou et al., 2023a). However, as these models become increasingly powerful, many have raised concerns that they could be repurposed by malicious actors to cause harm, motivating research on how to safeguard these models against malicious use.

Existing open-weight models often adapt safeguards designed for closed-weight models served through APIs (Touvron et al., 2023). These safeguards include refusal mechanisms and preference-based training, and they have provided substantial robustness against *input-based* jailbreaking attacks. However, recent work has demonstrated these safeguards are trivially defeated by attacks that edit model *weights*, breaking down after only a handful of fine-tuning steps (Qi et al., 2023). This poses a serious problem for open-weight models, because adversaries have full access to model weights and can tamper with built-in safeguards.

In this work, we study the problem of tamper-resistant safeguards for LLMs. This problem is depicted in Figure 1. Unlike existing research on LLM safeguards, we focus on attacks that modify model weights, which we refer to as tampering attacks. This problem has been considered very challenging and by some intractable, as no method has yet provided substantial robustness to these attacks. However, making progress on this problem would provide a valuable tool to regulators and model developers by ameliorating the dual-use dilemma of open-weight models (Miller & Selgelid, 2007).

To demonstrate that progress on this problem is possible, we develop the first LLM safeguards that obtain strong robustness against a wide variety of tampering attacks. Our approach allows developers to add a safeguard such that tampering attacks cannot easily remove the safeguard, while preserving the general capabilities of the LLM. We achieve this by performing adversarial training against tampering attacks, leveraging approaches from meta-learning. We identify various crucial factors that enable our method to work, including the choice of tamper-resistance loss, the selection of train-time adversaries, and the two-stage approach that we use for building in safeguards.
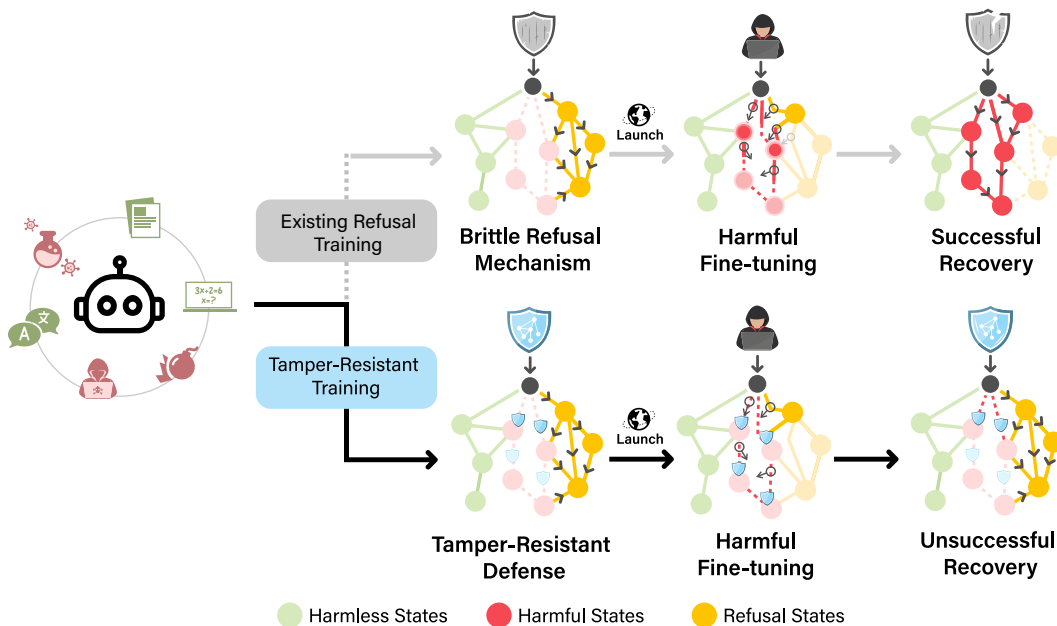
Figure 1: An illustration of existing brittle LLM safeguards compared to tamper-resistant training. Existing safeguards can be easily removed by adversaries with access to model weights, raising serious concerns about the safety of open-weight LLMs. We propose TAR, a method for making LLM safeguards robust to tampering attacks that fine-tune model weights.

We apply our method to develop tamper-resistant unlearning and refusal safeguards. In experiments, we demonstrate that our safeguards are far more robust to tampering attacks than prior methods. We stress-test our safeguards with extensive red teaming evaluations against 28 test-time adversaries, demonstrating resistance to fine-tuning attacks up to 5,000 steps. We hope our results foster future work on this important problem. Our experiment code and models are available at [anonymized].

## 2 RELATED WORK

**Adversarial attacks on LLMs.** Due to the extensive pre-training distribution of modern LLMs, they are prone to generating harmful content (Sheng et al., 2019; McGuffie & Newhouse, 2020). To mitigate this, many LLMs undergo fine-tuning to implement safeguards (Touvron et al., 2023; Bai et al., 2022; OpenAI, 2023), using methods such as reinforcement learning from human feedback (RLHF) (Christiano et al., 2017; Ouyang et al., 2022) and direct preference optimization (DPO) (Rafailov et al., 2023). While effective for normal use, these safeguards have been shown to be brittle, breaking down under jailbreak attacks (Wei et al., 2023; Zou et al., 2023b; Jin et al., 2024a) or a handful of fine-tuning steps on "uncensored" data (Qi et al., 2023; Zhan et al., 2023; Yang et al., 2023). This suggests current techniques for LLM alignment are inadequate, raising security concerns after deployment.

**LLM safeguards.** Since the discovery of these attacks, many safeguards have been proposed to defend against them. Against jailbreak attacks, defenses include system-level defenses Inan et al. (2023); Zhou et al. (2024); Yuan et al. (2024); Jain et al. (2023); Robey et al. (2023) that modify or filter model inputs or outputs and model-level defenses such as adversarial training Mazeika et al. (2024). Alternatively, some works explore machine unlearning as a way to remove harmful knowledge entirely with techniques such as influence functions (Koh & Liang, 2017; Bae et al., 2022), maximizing loss on forget sets (Yu et al., 2023; Eldan & Russinovich, 2023; Yao et al., 2023), or modifying representations (Wu et al., 2023; Belrose et al., 2024; Li et al., 2024; Sheshadri et al., 2024). However, jailbreaking defenses are not fully robust to adaptive adversaries Liu et al. (2023); Jin et al. (2024b), and existing unlearning methods are not robust to adversaries with access to model weights Lynch et al. (2024).

**Robust safeguards.** Several works have explored the tamper-resistance of unlearning methods for image classification (Golatkar et al., 2020a;b; Tarun et al., 2023). For bidirectional BERT-style models, Henderson et al. (2023) proposed a meta-learning approach for robustly preventing models from learning harmful tasks. In concurrent work, Deng et al. (2024) proposed a method extending this approach to small-scale vision classifiers and diffusion models. Recently, Liu et al. (2024) discussed the potential for robust unlearning in LLMs to improve the safety of open-source models, and Lynch et al. (2024) proposed evaluation metrics for robust unlearning in LLMs. To the best of our knowledge, no methods have been proposed for autoregressive LLMs that are robust to tampering attacks.

Several concurrent works have explored ways of defending LLM refusal mechanisms against fine-tuning. Huang et al. (2024) add a perturbation loss to make an LLM learn to produce embeddings that are more invariant to perturbations, Rosati et al. (2024a) maximize prediction loss on harmful generations while minimizing loss on refusals, and Rosati et al. (2024b) regularize harmful representations to look random. Unfortunately, these works evaluate against small sets of fine-tuning adversaries or have limited robustness. We corroborate this in our comparisons, finding that the approaches in the latter two works lack robustness to the tampering attacks in our evaluations.



Figure 2: Comparison of our TAR method to 12 baseline safeguards. Unlike prior methods, TAR provides far greater tamper-resistance at similar levels of general capability, measured via MMLU. Tamper-resistance is computed as the normalized error on WMDP Biosecurity, Chemical Security, and Cybersecurity questions (Li et al., 2024), averaged across up to 28 fine-tuning attacks.

## 3 TAMPER-RESISTANT SAFEGUARDS

### 3.1 THREAT MODEL

We assume the defender releases an LLM with weights $\theta_G$ and a safeguard $G$ applied. The defender's goal is to design $G$ such that $\theta_G$ obtains high values on `safety_metric`$(\theta_G)$ and `capabilities_metric`$(\theta_G)$. Moreover, the defender seeks to preserve a high value of `safety_metric`$(\theta_G)$ after the adversary's move. We consider a compute-bounded adversary with unrestricted access to $\theta_G$, enabling attacks that directly modify $\theta_G$. We refer to these as "tampering attacks." The adversary's goal is to obtain a model $\theta'_G$ that minimizes the safety metric given reasonable compute limits, such as fine-tuning for 1,000 to 5,000 steps. We assume the adversary will not spend a significant fraction of the compute required to pre-train the LLM, since at that point they could train their own model without safeguards.

### 3.2 PROBLEM DEFINITION AND METRICS

We describe a general notation for quantifying the tamper-resistance of safeguards. Define $G$, $\theta_G$, `safety_metric`, and `capabilities_metric` as in the threat model. Let `attack` denote a compute-bounded adversarial attack that maps $\theta_G$ to $\theta'_G$, with stronger attacks obtaining lower values of `safety_metric`$(\theta'_G)$. We say that a safeguard $G$ is *tamper-resistant* if its post-attack `safety_metric`$(\theta'_G)$ is high across a broad range of strong test-time adversarial attacks $\mathcal{A}_{\text{test}}$.

Note that $\theta_G$ often modifies an underlying $\theta$ that lacks safeguards, often through a fine-tuning procedure. Additionally, strong tamper-resistance can be obtained if the safeguard simply overwrites $\theta$ with noise, but this model would no longer be useful. Thus, maintaining a high `capabilities_metric`$(\theta_G)$ is crucial, and evaluation of a safeguard must consider both its tamper-resistance and how well it preserves general capabilities.

We focus on two common safeguard domains: weaponization knowledge restriction and harmful request refusal. In each domain, we define safety and capabilities test metrics, which we use alongside test-time adversaries to evaluate tamper-resistant safeguards.
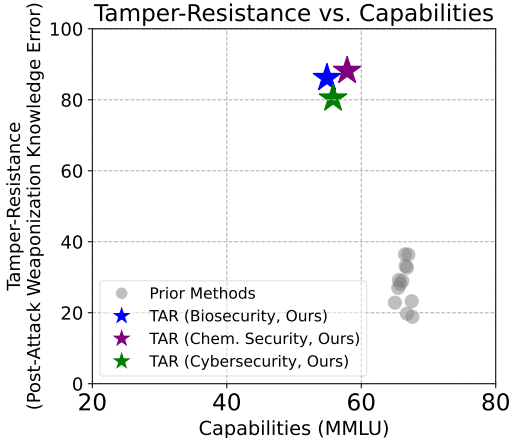
**Weaponization knowledge restriction.** In weaponization knowledge restriction, safeguards prevent the model from producing text about weaponization knowledge, while preserving capabilities for benign knowledge domains. Existing safeguards of this nature include representation engineering methods like circuit breaking (Zou et al., 2024). The `safety_metric` is defined as error on a *forget set*, and the `capabilities_metric` is defined as accuracy on a *retain set*. Specifically, we consider the problem of restricting biosecurity, chemical security, and cybersecurity knowledge, and evaluate the resulting model on the *Weapons of Mass Destruction Proxy (WMDP)* benchmark (Li et al., 2024). *WMDP* contains 3,668 multiple-choice questions, spanning biosecurity, chemical security, and cybersecurity knowledge. Importantly, *WMDP* questions do not evaluate hazardous knowledge directly, but instead measure proxy expert-level knowledge for each hazardous domain, such that restricting the expert-level knowledge would also restrict the hazardous knowledge. We define the forget set as the respective hazardous knowledge subject in *WMDP*, and retain set as the complement of the given subject in *MMLU* (Hendrycks et al., 2021), a multi-task question-answering benchmark spanning 57 tasks across a variety of knowledge domains.

**Harmful request refusal.** In the harmful request refusal setting, safeguards prevent the model from producing "harmful" outputs. We define the `safety_metric` as the complement of average Attack Success Rate (ASR) of various jailbreaking attacks, while the `capabilities_metric` captures the conversational abilities of $\theta_G$. Specifically, we use a static set of test cases from *HarmBench*, an automated red-teaming framework for measuring prompt jailbreak robustness in LLMs, to evaluate jailbreak ASR (Mazeika et al., 2024) after tampering attacks. We use *MT-Bench*, a multi-turn question-answering benchmark graded by an LLM judge, to evaluate conversational abilities (Zheng et al., 2023a).

### 3.3 RED TEAMING

To properly measure the robustness of tamper-resistant safeguards, we conduct red-teaming with up to 28 adversaries, including many that are unseen at training time. In our evaluations, we subject our method to adversaries with varying compute budgets, access to held-out datasets, and diverse hyperparameters. For fine-tuning adversaries, we vary the learning rate, learning rate scheduler, optimization algorithm, and batch size. Many of these adversaries were fixed during early experiments, with some added over time as we found attacks that broke intermediate versions of our method. Extensive stress testing of this nature is critical for obtaining confidence in a tamper-resistant safeguard. For research on developing these safeguards, extensive red teaming also allows measuring incremental progress, using the number and strength of existing attacks one can defend against as a robustness metric.

## 4 SAFEGUARD TAMPER-RESISTANCE TRAINING

To obtain tamper-resistant safeguards, we propose a new method outlined in Algorithm 1 inspired by adversarial training and meta-learning to directly strengthen LLM safeguards against tampering attacks, called Tampering Attack Resistance (TAR). We identify unique properties of this adversarial training regime and leverage them to improve robustness.

Our method for training tamper-resistant safeguards consists of two phases: (1) model safeguarding and (2) tamper-resistance training.

### 4.1 MODEL SAFEGUARDING

The method begins by including an initial safeguard $G$ into a base model $\theta$. For example, initial safeguards for knowledge restriction can be drawn from a wide variety of existing methods, including circuit breaking (Li et al., 2024) or constrained gradient ascent for a particular knowledge domain. Similarly, we can include a refusal safeguard by performing RLHF (Ouyang et al., 2022) or DPO (Rafailov et al., 2023) on refusal completions. Importantly, these initial safeguards do not need to be tamper-resistant. Empirically, we find that the safeguarding step is important for strong tamper-resistance.

---

**Algorithm 1** TAR: Tampering Attack Resistance

---

**Input:** Initial LLM parameters $\theta$, train-time adversary set $\mathcal{A}_{\text{train}}$, `capabilities_metric` proxy dataset $\mathcal{D}_{\text{retain}}$, `safety_metric` proxy dataset $\mathcal{D}_{\text{TR}}$, outer steps $N$, learning rate $\eta$, number of sampled adversaries $K$, tamper-resistance loss scale $\lambda_{\text{TR}}$, retain loss scale $\lambda_{\text{retain}}$, $h_{\theta}(\cdot)$ returns the residual stream hidden states for model parameters $\theta$

$\theta_0 \leftarrow$ **Apply Initial Safeguard to** $\theta$
**for** $i = 1$ **to** $N$ **do**
    $g_{\text{TR}} \leftarrow 0$ *# For accumulating tamper-resistance gradient*
    Sample $x_{\text{TR}} \sim \mathcal{D}_{\text{TR}}$
    **for** $k = 1$ **to** $K$ **do**
        Sample `attack` $\sim \mathcal{A}_{\text{train}}$
        *# Tamper-resistance loss from Equation 1*
        $g_{\text{TR}} \leftarrow g_{\text{TR}} + \frac{1}{K} \nabla_{\theta_{i-1}} \mathcal{L}_{\text{TR}}(\texttt{attack}(\theta_{i-1}), x_{\text{TR}})$
    **end for**
    Sample $x_r \sim \mathcal{D}_{\text{retain}}$
    *# RepE retain loss from Equation 2*
    $g_{\text{retain}} \leftarrow \nabla_{\theta_{i-1}} \Big( \mathcal{L}_{\text{LM}}(\theta_{i-1}, x_r) + \|h_{\theta_{i-1}}(x_r) - h_{\theta}(x_r)\|_2^2 \Big)$
    *# Full tamper-resistance update*
    $\theta_i \leftarrow \theta_{i-1} - \eta \Big( \lambda_{\text{TR}} \cdot g_{\text{TR}} + \lambda_{\text{retain}} \cdot g_{\text{retain}} \Big)$
**end for**
$\theta_G \leftarrow \theta_N$
**return** $\theta_G$

---

## 4.2 TAMPER-RESISTANCE TRAINING

Starting from $\theta_{G_0}$, we train the tamper-resistant $\theta_G$ using a novel adversarial training procedure. Namely, we train against a set of tampering attacks $\mathcal{A}_{\text{train}}$, where the defender's objective is to maximize a proxy `safety_metric` after applying an adversarial attack `attack` $\sim \mathcal{A}_{\text{train}}$ to $\theta$. Since it may not be feasible to differentiate through `attack`, we draw on insights from prior work in meta-learning, defining $\texttt{attack}(\theta_G) = \theta'_G = \theta_G + \texttt{attack}'(\theta_G)$ as a perturbation on top of initial parameters, where backpropagation through $\texttt{attack}'$ is approximated with a straight-through estimator (Bengio et al., 2013).

We focus on supervised fine-tuning (SFT) adversaries where `attack` applies several steps of optimization to $\theta_G$, which allows straight-through estimation through $\texttt{attack}'$ to benefit from the setting and approximations of first-order MAML (Finn et al., 2017). However, we note key differences in our approach from standard meta-learning and prior methods (Finn et al., 2017; Henderson et al., 2023). In particular, traditional meta-learning techniques seek to obtain a model initialization that is *close* to optimality on multiple test distributions. In our setting, we seek to obtain an initialization that is *far* from optimality on multiple adversaries' test distributions. Novel to our approach in this new setting is the use of a tamper-resistance loss in the "outer loop" that differs from the fine-tuning adversary's loss function and serves to maximize the proxy safety metric. We depict this structure in Algorithm 1, and explain the objective below.

**Impeding the adversary's loss.** The aim of tamper-resistance training is to prevent adversaries with large compute budgets from reducing the `safety_metric` at *test-time*. In adversarial training for tamper-resistance, we define a tamper-resistance loss $\mathcal{L}_{\text{TR}}$ that counters `attack`. We operationalize our goal of *avoiding adversary optimality* as searching for $\theta$ such that $\mathcal{L}_{\text{TR}}$ is minimized for $\texttt{attack}(\theta)$.

Empirically, we find that the choice of tamper-resistance loss $\mathcal{L}_{\text{TR}}$ significantly affects this goal. Prior work (Henderson et al., 2023) negates the loss of a fine-tuning adversary, in which the aim is to arbitrarily maximize the adversary's loss throughout fine-tuning. This formulation has two issues: (1) maximizing a cross-entropy loss can cause divergence; (2) empirically we observe that when using this objective against fine-tuning adversaries, the model learns to explode the adversary's loss for the first few inner loop steps, while loss at later steps remains low. In Figure 3, we show the difference in choosing $\mathcal{L}_{\text{TR}}$ to be a clamped negative cross-entropy loss vs. negative entropy loss for weaponization
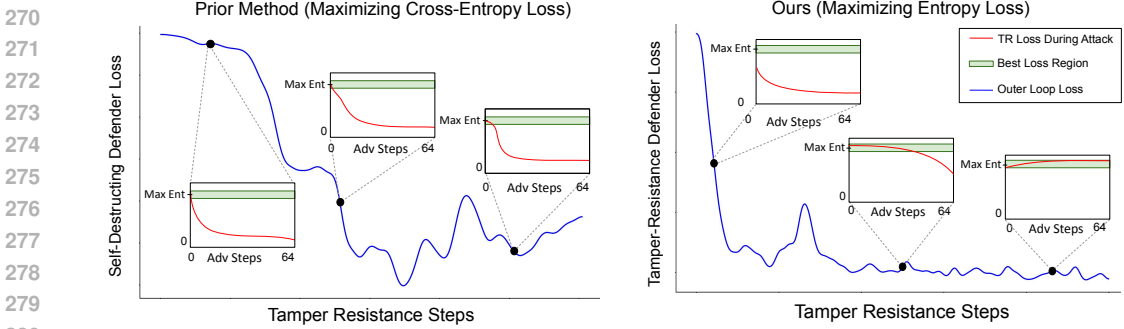
Figure 3: The choice of tamper-resistance loss is crucial for obtaining good performance. Here, we show loss trajectories when the tamper-resistance loss is negative cross-entropy (left), versus negative entropy (right), over the course of TAR for 750 steps. Outer loop losses (blue) are reduced by the defender, and inner-loop losses (red) are reduced by the train-time adversary. When the tamper-resistance loss maximizes cross-entropy (left), the adversary is only affected earlier in its trajectory and quickly recovers. By contrast, when the tamper-resistance loss maximizes entropy (right), the inner loop adversary is eventually thwarted along its entire trajectory. Plots are smoothed.

knowledge restriction. For the latter, $\mathcal{L}_{\text{TR}}$ is eventually satisfied for all inner loop steps. For harmful request refusal, we choose $\mathcal{L}_{\text{TR}}$ to be the DPO loss (Rafailov et al., 2023). We provide further detail on the choice of $\mathcal{L}_{\text{TR}}$ in both settings in Appendix B.2.

**Tamper-resistance objective.** We now describe the general proxy objective used for preventing a tampering attack from recovering weaponization knowledge or harmful behavior. For a given `safety_metric`, let $\mathcal{D}_{\text{TR}}$ and $\mathcal{L}_{\text{TR}}$ respectively be a dataset and loss function such that minimizing $\mathcal{L}_{\text{TR}}(\theta_G, \mathcal{D}_{\text{TR}})$ serves as a proxy objective for maximizing the `safety_metric`$(\theta_G)$. We define $\mathcal{D}_{\text{retain}}$ and $\mathcal{L}_{\text{retain}}$ correspondingly for `capabilities_metric`$(\theta_G)$. The defender's objective is to solve the following optimization problem:

$$\min_{\theta} \lambda_{\text{TR}} \cdot \mathbb{E}_{\text{attack} \sim \mathcal{A}_{\text{train}}} \big[ \mathcal{L}_{\text{TR}}(\texttt{attack}(\theta); \mathcal{D}_{\text{TR}}) \big] + \lambda_{\text{retain}} \cdot \mathcal{L}_{\text{retain}}(\theta; \mathcal{D}_{\text{retain}}), \qquad (1)$$

where $\mathcal{L}_{\text{TR}}$ is a tamper-resistance loss that counters $\texttt{attack}(\theta)$. The $\mathcal{L}_{\text{retain}}$ term is a representation engineering (Zou et al., 2023a) inspired retain loss for preserving performance on the capabilities proxy dataset $\mathcal{D}_{\text{retain}}$, given by

$$\mathcal{L}_{\text{retain}}(\theta; \mathcal{D}_{\text{retain}}) = \mathbb{E}_{x \sim \mathcal{D}_{\text{retain}}} \Big[ \mathcal{L}_{\text{LM}}(\theta, x) + \| h_\theta(x) - h_{\theta_{G_0}}(x) \|_2^2 \Big] \qquad (2)$$

where $h_\theta(\cdot)$ returns the residual stream hidden states for model parameters $\theta$ and $\mathcal{L}_{\text{LM}}$ is the standard language modelling cross-entropy loss. Empirically, we find that pushing retain-set residual stream representations to be close to the base model $\theta_{G_0}$ via the $\ell_2$-norm loss in Equation 2 maintains a high `capabilities_metric`$(\theta_G)$. In Equation 1 we include $\lambda_{\text{TR}}$ and $\lambda_{\text{retain}}$ as scalar weightings for the tamper-resistance loss and retain loss, respectively. We provide further details on the design of the tamper-resistance loss function in Appendix B.2 as well as an efficiency trick for sampling fine-tuning attacks for TAR in Appendix B.3.

## 5  EXPERIMENTS

We evaluate TAR in weaponization knowledge restriction and harmful request refusal settings, with results shown in Table 1 and Table 2 respectively. We discuss the setup, baselines, and analysis for our results. In each setting, we use a specific set of training adversaries $\mathcal{A}_{\text{train}}$ and test adversaries $\mathcal{A}_{\text{test}}$. Further experiment details are presented in Appendix D.

### 5.1  WEAPONIZATION KNOWLEDGE RESTRICTION

We now describe the setup, baselines, and results for our weaponization knowledge restriction experiments, including the knowledge domains, optimizers, and evaluation details.

| Weaponization Domain | Model | Pre-Attacks | | Post-Attacks (Avg) |
|---|---|---|---|---|
| | | Retain ($\uparrow$) | Forget ($\downarrow$) | Forget ($\downarrow$) |
| Biosecurity | Random | 25.0 | 25.0 | 25.0 |
| | No Defense | 67.3 | 70.5 | 70.5 |
| | Max Entropy | 65.0 | 33.2 | 60.1 |
| | Min Posterior | 65.6 | 50.4 | 57.2 |
| | LLMU | 65.5 | <u>29.9</u> | 58.2 |
| | RMU | 65.8 | 31.2 | 57.7 |
| | TAR (Ours) | 54.9 | **24.0** | **31.3** |
| Chemical Security | Random | 25.0 | 25.0 | 25.0 |
| | No Defense | 68.2 | 47.8 | 47.8 |
| | Max Entropy | 67.5 | 50.0 | 43.1 |
| | Min Posterior | 66.8 | 49.5 | 43.7 |
| | LLMU | 67.0 | 30.1 | 42.3 |
| | RMU | 67.6 | **27.5** | 44.6 |
| | TAR (Ours) | 57.9 | <u>29.2</u> | **27.6** |
| Cybersecurity | Random | 25.0 | 25.0 | 25.0 |
| | No Defense | 68.2 | 46.4 | 46.4 |
| | Max Entropy | 66.5 | 28.7 | 39.4 |
| | Min Posterior | 66.6 | 41.8 | 39.7 |
| | LLMU | 66.1 | <u>27.6</u> | 40.6 |
| | RMU | 66.8 | 29.5 | 39.9 |
| | TAR (Ours) | 55.8 | **26.3** | **29.5** |

Table 1: Pre-Attack and average Post-Attack accuracies for WMDP Biosecurity, Chemical Security, and Cybersecurity for TAR and all other baselines, reported for Llama-3-8B. The average Post-Attack accuracy is computed as the average accuracy across the 28 fine-tuning attacks discussed in Section 5. TAR is the only method that maintains low Post-Attack recovery while preserving high Retain MMLU and low Forget accuracies. All values are percentages.

**Setup.** We focus on implementing tamper-resistant safeguards for restricting proxy weaponization knowledge about biosecurity, chemical security, and cybersecurity from Llama-3-8B-Instruct (AI@Meta, 2024) that has been initially safeguarded via the *Random Mapping* method discussed in Appendix B.1. For each weaponization domain, we assign $\mathcal{D}_{TR}$ to the corresponding forget set described in Appendix D.1. We proceed to sample train-time 64-step fine-tuning attacks from different data distributions, detailed in Appendix D.2. We use $N = 750$ outer loop steps, ScheduleFree AdamW (Defazio et al., 2024) with a learning rate of $2 \times 10^{-5}$ as the outer loop tamper-resistance optimizer, and loss scales of $\lambda_{TR} = 4.0$, $\lambda_{\text{retain}} = 1.0$. Lastly, We evaluate Pre-Attack and Post-Attack accuracy on corresponding WMDP subjects (Li et al., 2024) averaged across all adversaries in Appendix E.1, and measure benign capabilities via the complement of subjects related to each proxy weaponization domain in MMLU (Hendrycks et al., 2021).

**Baselines.** We evaluate two recently proposed knowledge restriction methods: *RMU* (Li et al., 2024) and *LLMU* (Yao et al., 2023). We also design two baseline methods for knowledge restriction: *Min Posterior*, which minimizes posterior loss on forget set tokens; *Max Entropy*, which maximizes entropy on forget set tokens. Two additional methods, *MLAC* (Henderson et al., 2023) and *SOPHON* (Deng et al., 2024), require substantial modifications for the LLM setting, so we show results on adapted versions of these baselines in Appendix F.3.

**Results.** We show weaponization knowledge restriction safeguard results on Llama-3-8B-Instruct in Table 1 and Figure 2. These results are averaged across all 28 adversaries from Appendix E.1. Our large-scale experiments corroborate the findings in recent work that existing LLM safeguards are extremely brittle to fine-tuning attacks. By contrast, TAR maintains low post-attack forget accuracy across all three domains. However, we observe that TAR lowers retain accuracy by 11.3% on average, indicating a trade-off between benign capabilities and robustness. In Figure 4, we observe that TAR is robust to significantly more fine-tuning attacks than all prior methods. While existing baselines
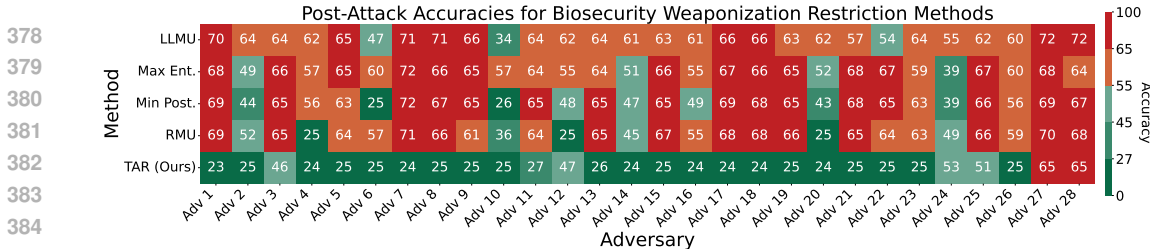
Figure 4: Red teaming results for weaponization knowledge restriction in the biosecurity domain. All numbers are percentages, with ideal defenses obtaining 25%. Red colors indicate that an attack recovered performance near the level of the No Defense baseline. We evaluate each defense against a diverse range of strong adversaries, described in Appendix E.1. Compared to prior safeguards, TAR greatly increases tamper-resistance for nearly all adversaries.

break down under most attacks, TAR obtains a post-attack forget accuracy near random chance for nearly all attacks, indicating a successful defense.

Overall, we find that TAR provides significantly more robustness to realistic fine-tuning attacks than all prior methods, including attacks as strong as 5,000 steps of optimization on completely held-out data. These results demonstrate for the first time that obtaining strong tamper-resistance for open-weight LLMs may be possible.

## 5.2 HARMFUL REQUEST REFUSAL

|  | Refusal Trained | R2D2 | RepNoise | RR | TAR (Ours) |
|---|---|---|---|---|---|
| Pre-Attacks MT-Bench ($\uparrow$) | 8.1 | 6.0 | 6.2 | 8.0 | 6.3 |
| Avg. Post-Attacks ASR ($\downarrow$) | 72.5 | 78.3 | 74.5 | 84.8 | **63.9** |

Table 2: Average Post-Attack HarmBench ASR, reported for TAR, Representation Rerouting (RR), and the Refusal Trained Llama-3-8B-Instruct model across 5 fine-tuning attacks depicted in Appendix E.2, as well as Pre-Attack MT-Bench. TAR is more robust than other methods after tampering, while maintaining comparable MT-Bench performance. ASR values are percentages.

We now describe the setup, baselines, and results for our harmful request refusal experiments, including the datasets used and evaluation details.

**Setup.** For harmful request refusal training, we seek to make existing refusal safeguards in Llama-3-8B-Instruct robust to tampering attacks. We sample train-time adversaries that perform 64-step SFT attacks using the Anthropic-HH-RLHF dataset (Bai et al., 2022), following the methodology in Appendix D.2. Similar to the weaponization knowledge restriction setting, we use $N = 100$ outer loop steps, ScheduleFree AdamW (Defazio et al., 2024) with an LR of $6 \times 10^{-5}$ as the outer loop tamper-resistance optimizer, and loss scales of $\lambda_{\text{TR}} = 0.1, \lambda_{\text{retain}} = 1.0$. We evaluate the Post-Attack jailbreak attack success rate (ASR) on HarmBench (Mazeika et al., 2024) after the tampering attacks in Appendix E.2, and measure benign capabilities preservation via MT-Bench (Zheng et al., 2023b), which evaluates multi-turn conversation ability.

**Baselines.** We consider 4 baselines alongside our TAR model: Llama-3-8B-Instruct *(Refusal Trained)*; *Representation Rerouting (RR)* (Zou et al., 2024) on Llama-3-8B-Instruct, which trains to push representations for harmful input prompts to be orthogonal to the original corresponding representations in Llama-3-8B-Instruct; *R2D2* (Mazeika et al., 2024) on Zephyr-7B (Tunstall et al., 2023), which performs adversarial training against GCG attacks (Zou et al., 2023b); and *RepNoise* (Rosati et al., 2024b) on Llama-2-7B (Touvron et al., 2023), which regularizes harmful representations to Gaussian noise.

**Results.** We show refusal results in Table 2. While the Refusal Training, RR, and R2D2 baselines resist jailbreak attacks in HarmBench before tampering, we find that percentage attack success

rate jumps up to above 77 after tampering, while our TAR method only rises to 61.7. Since we apply our TAR refusal safeguard to Llama-3-8B, it does reduce MT-Bench by 1.7. However, this exceeds the MT-Bench score of fairly capable open-weight models, indicating that benign capabilities are largely preserved. We leave the exploration of the full impact on capabilities to future work. Additional results are in Table 11. In general, we find that our TAR model refuses more Post-Attack jailbreaks than previous baselines, and demonstrates the flexibility of the tamper-resistance objective to accommodate the harmful request refusal setting.

## 5.3 Analysis

**Red teaming.** To stress-test the tamper-resistance of our models, we conduct an extensive suite of supervised fine-tuning attacks with 28 distinct adversaries. We vary the optimizer, number of optimization steps, learning rate, learning rate schedule, fine-tuning dataset, batch size, and overall fine-tuning method (e.g., full fine-tuning versus parameter-efficient fine-tuning). By default, our attacks use $1,000$ fine-tuning steps, although some use $5,000$ steps. Full details for these adversaries are provided in Table 9.



Figure 5: Our TAR safeguard remains robust to fine-tuning attacks that greatly exceed the $64$ steps used by train-time adversaries. For the LLMU safeguard, the adversary's loss quickly decreases into the recovery region. By contrast, TAR maintains flatness at a high loss for all $1,000$ steps.

We show red teaming results in Figures 4 and 8. While baseline safeguards withstand fine-tuning attacks in a small number of cases, most adversaries succeed in removing the safeguards. By contrast, our TAR safeguard is robust to a wide range of adversaries. This shows that tamper-resistance is a tractable problem on which progress can be made. However, we find that when sampling solely full-parameter SFT adversaries during TAR, the resulting model is not robust to parameter-efficient fine-tuning (PEFT) adversaries (adversaries 27 and 28), highlighting the importance of extensive red teaming when developing tamper-resistant defenses. We hypothesize that future work can easily address this issue, as we demonstrate in Appendix C.3 that targeted patching of vulnerabilities is possible.

**Generalization to stronger test-time attacks.** In Figure 5, we show a fine-tuning attack at an LR of $2 \times 10^{-5}$ that attempts to recover biosecurity knowledge on our TAR model and a model safeguarded with LLMU. We find that the tamper-resistance of TAR generalizes far beyond the $64$ steps used by train-time adversaries. Surprisingly, we observe that the test-time adversary's cross-entropy loss does not decrease below 7 for all $1,000$ steps. Moreover, the loss enters a plateau and does not decrease at all after 200 steps. In Appendix C.2, we find that the length and height of this plateau can be increased by increasing the number of inner loop steps during adversarial training. In contrast, with LLMU as the safeguard, the adversary's loss decreases to within the recovery region in under 20 steps. We note that the adversary's test loss decreasing into the recovery region does not always correspond with recovery on downstream metrics (e.g., WMDP), but often does.

## 6 Conclusion

We introduced a novel method for implementing tamper-resistant safeguards for LLMs and explored applications in weaponization knowledge restriction and harmful refusal training. We compare our results to prior work in each setting, finding that our method is the first method robust under the rigorous red-teaming evaluation that we consider. More broadly, we demonstrate that progress on open-weight tamper-resistance is tractable. We believe this line of research is crucial for enabling ongoing deployment of robust, open-weight LLMs, ensuring their alignment with regulatory frameworks and preemptively addressing the risk of malicious use.
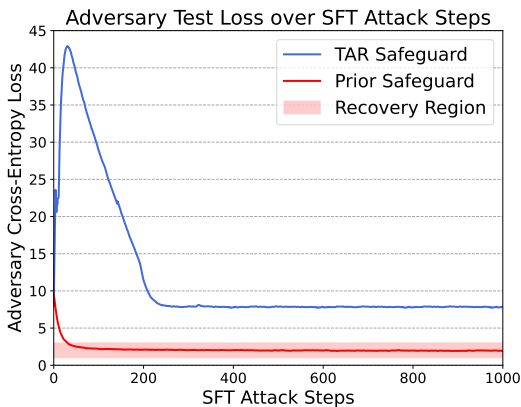
# REFERENCES

AI@Meta. Llama 3 model card. 2024.

Juhan Bae, Nathan Ng, Alston Lo, Marzyeh Ghassemi, and Roger Baker Grosse. If influence functions are the answer, then what is the question? *ArXiv*, abs/2209.05364, 2022.

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

Nora Belrose, David Schneider-Joseph, Shauli Ravfogel, Ryan Cotterell, Edward Raff, and Stella Biderman. Leace: Perfect linear concept erasure in closed form. *Advances in Neural Information Processing Systems*, 36, 2024.

Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation, 2013.

Tianchi Cai, Xierui Song, Jiyan Jiang, Fei Teng, Jinjie Gu, and Guannan Zhang. Ulma: Unified language model alignment with human demonstration and point-wise preference, 2024.

Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.

CTFtime. Ctftime writeups archive. https://ctftime.org/writeups, 2024.

Aaron Defazio, Xingyu, Yang, Harsh Mehta, Konstantin Mishchenko, Ahmed Khaled, and Ashok Cutkosky. The road less scheduled, 2024.

Jiangyi Deng, Shengyuan Pang, Yanjiao Chen, Liangming Xia, Yijie Bai, Haiqin Weng, and Wenyuan Xu. Sophon: Non-fine-tunable learning to restrain task transferability for pre-trained models, 2024.

Ronen Eldan and Mark Russinovich. Who's harry potter? approximate unlearning in llms. *ArXiv*, abs/2310.02238, 2023.

Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks, 2017.

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The pile: An 800gb dataset of diverse text for language modeling, 2020.

Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9304–9312, 2020a.

Aditya Golatkar, Alessandro Achille, and Stefano Soatto. Forgetting outside the box: Scrubbing deep networks of information accessible from input-output observations. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIX 16*, pp. 383–398. Springer, 2020b.

Peter Henderson, Eric Mitchell, Christopher D. Manning, Dan Jurafsky, and Chelsea Finn. Self-destructing models: Increasing the costs of harmful dual uses of foundation models, 2023.

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

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021.

Gao Huang, Yixuan Li, Geoff Pleiss, Zhuang Liu, John E. Hopcroft, and Kilian Q. Weinberger. Snapshot ensembles: Train 1, get m for free, 2017.

Tiansheng Huang, Sihao Hu, and Ling Liu. Vaccine: Perturbation-aware alignment for large language model, 2024.

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

Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping yeh Chiang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. Baseline defenses for adversarial attacks against aligned language models, 2023.

Haibo Jin, Ruoxi Chen, Andy Zhou, Yang Zhang, and Haohan Wang. Guard: Role-playing to generate natural-language jailbreakings to test guideline adherence of large language models, 2024a.

Haibo Jin, Andy Zhou, Joe D. Menke, and Haohan Wang. Jailbreaking large language models against moderation guardrails via cipher characters, 2024b.

Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.

Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International Conference on Machine Learning*, 2017.

Solomon Kullback and R. A. Leibler. On information and sufficiency. *Annals of Mathematical Statistics*, 22:79–86, 1951.

Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for "mind" exploration of large language model society, 2023.

Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D. Li, Ann-Kathrin Dombrowski, Shashwat Goel, Long Phan, Gabriel Mukobi, Nathan Helm-Burger, Rassin Lababidi, Lennart Justen, Andrew B. Liu, Michael Chen, Isabelle Barrass, Oliver Zhang, Xiaoyuan Zhu, Rishub Tamirisa, Bhrugu Bharathi, Adam Khoja, Zhenqi Zhao, Ariel Herbert-Voss, Cort B. Breuer, Andy Zou, Mantas Mazeika, Zifan Wang, Palash Oswal, Weiran Liu, Adam A. Hunt, Justin Tienken-Harder, Kevin Y. Shih, Kemper Talley, John Guan, Russell Kaplan, Ian Steneker, David Campbell, Brad Jokubaitis, Alex Levinson, Jean Wang, William Qian, Kallol Krishna Karmakar, Steven Basart, Stephen Fitz, Mindy Levine, Ponnurangam Kumaraguru, Uday Tupakula, Vijay Varadharajan, Yan Shoshitaishvili, Jimmy Ba, Kevin M. Esvelt, Alexandr Wang, and Dan Hendrycks. The wmdp benchmark: Measuring and reducing malicious use with unlearning, 2024.

Sijia Liu, Yuanshun Yao, Jinghan Jia, Stephen Casper, Nathalie Baracaldo, Peter Hase, Xiaojun Xu, Yuguang Yao, Hang Li, Kush R Varshney, et al. Rethinking machine unlearning for large language models. *arXiv preprint arXiv:2402.08787*, 2024.

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

Llama Team, AI @ Meta. The llama 3 herd of models, 2024.

Ilya Loshchilov and Frank Hutter. SGDR: stochastic gradient descent with restarts. *CoRR*, abs/1608.03983, 2016.

Aengus Lynch, Phillip Guo, Aidan Ewart, Stephen Casper, and Dylan Hadfield-Menell. Eight methods to evaluate robust unlearning in llms. *arXiv preprint arXiv:2402.16835*, 2024.

Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, David Forsyth, and Dan Hendrycks. Harmbench: A standardized evaluation framework for automated red teaming and robust refusal. In *ICML*, 2024.

Kris McGuffie and Alex Newhouse. The radicalization risks of gpt-3 and advanced neural language models, 2020.

Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models, 2016.

Seumas Miller and Michael J Selgelid. Ethical and philosophical consideration of the dual-use dilemma in the biological sciences. *Science and engineering ethics*, 13:523–580, 2007.

Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms, 2018.

OpenAI. Gpt-4 technical report, 2023.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.

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

Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In *Neural Information Processing Systems (NeurIPS)*, 2023.

Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations toward training trillion parameter models. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 1–16, 2020. doi: 10.1109/SC41405.2020.00024.

Jie Ren, Samyam Rajbhandari, Reza Yazdani Aminabadi, Olatunji Ruwase, Shuangyan Yang, Minjia Zhang, Dong Li, and Yuxiong He. ZeRO-Offload: Democratizing Billion-Scale model training. In *2021 USENIX Annual Technical Conference (USENIX ATC 21)*, pp. 551–564. USENIX Association, July 2021. ISBN 978-1-939133-23-6.

Alexander Robey, Eric Wong, Hamed Hassani, and George J. Pappas. Smoothllm: Defending large language models against jailbreaking attacks, 2023.

Domenic Rosati, Jan Wehner, Kai Williams, Łukasz Bartoszcze, Jan Batzner, Hassan Sajjad, and Frank Rudzicz. Immunization against harmful fine-tuning attacks. *arXiv preprint arXiv:2402.16382*, 2024a.

Domenic Rosati, Jan Wehner, Kai Williams, Łukasz Bartoszcze, David Atanasov, Robie Gonzales, Subhabrata Majumdar, Carsten Maple, Hassan Sajjad, and Frank Rudzicz. Representation noising effectively prevents harmful fine-tuning on llms, 2024b.

Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. The woman worked as a babysitter: On biases in language generation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2019.

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

Ayush K Tarun, Vikram S Chundawat, Murari Mandal, and Mohan Kankanhalli. Fast yet effective machine unlearning. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv:2307.09288*, 2023.

Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. Zephyr: Direct distillation of lm alignment, 2023. URL https://arxiv.org/abs/2310.16944.

Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. Openchat: Advancing open-source language models with mixed-quality data, 2023.

Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does llm safety training fail? In *Neural Information Processing Systems (NeurIPS)*, 2023.

Xinwei Wu, Junzhuo Li, Minghui Xu, Weilong Dong, Shuangzhi Wu, Chao Bian, and Deyi Xiong. Depn: Detecting and editing privacy neurons in pretrained language models. In *Conference on Empirical Methods in Natural Language Processing*, 2023.

Xingyu Xie, Pan Zhou, Huan Li, Zhouchen Lin, and Shuicheng Yan. Adan: Adaptive nesterov momentum algorithm for faster optimizing deep models, 2023.

Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. Magpie: Alignment data synthesis from scratch by prompting aligned llms with nothing, 2024.

Xianjun Yang, Xiao Wang, Qi Zhang, Linda Petzold, William Yang Wang, Xun Zhao, and Dahua Lin. Shadow alignment: The ease of subverting safely-aligned language models, 2023. URL https://arxiv.org/abs/2310.02949.

Yuanshun Yao, Xiaojun Xu, and Yang Liu. Large language model unlearning. *ArXiv*, abs/2310.10683, 2023.

Charles Yu, Sullam Jeoung, Anish Kasi, Pengfei Yu, and Heng Ji. Unlearning bias in language models by partitioning gradients. In *Annual Meeting of the Association for Computational Linguistics*, 2023.

Zhuowen Yuan, Zidi Xiong, Yi Zeng, Ning Yu, Ruoxi Jia, Dawn Song, and Bo Li. Rigorllm: Resilient guardrails for large language models against undesired content, 2024.

Matthew D. Zeiler. ADADELTA: an adaptive learning rate method. *CoRR*, abs/1212.5701, 2012.

Qiusi Zhan, Richard Fang, Rohan Bindu, Akul Gupta, Tatsunori Hashimoto, and Daniel Kang. Removing rlhf protections in gpt-4 via fine-tuning, 2023.

Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright, Hamid Shojanazeri, Myle Ott, Sam Shleifer, Alban Desmaison, Can Balioglu, Pritam Damania, Bernard Nguyen, Geeta Chauhan, Yuchen Hao, Ajit Mathews, and Shen Li. Pytorch fsdp: Experiences on scaling fully sharded data parallel, 2023.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023a.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023b.

Andy Zhou, Bo Li, and Haohan Wang. Robust prompt optimization for defending language models against jailbreaking attacks, 2024.

Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, J. Zico Kolter, and Dan Hendrycks. Representation engineering: A top-down approach to ai transparency, 2023a.

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

Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, Rowan Wang, Zico Kolter, Matt Fredrikson, and Dan Hendrycks. Improving alignment and robustness with short circuiting. *arXiv preprint arXiv:2406.04313*, 2024.

13

## A LIMITATIONS

Our method for training tamper-resistant safeguards demonstrates considerable robustness against a wide range of tampering attacks, yet several avenues for improvement remain: (1) While we focus on supervised fine-tuning attacks, the broader spectrum of open-weight tampering techniques necessitates diverse future red-teaming efforts. (2) Scaling to larger models poses computational challenges that require optimization to reduce overheads.

Additionally, in cases where TAR maintains a low post-attack forget accuracy, the post-attack retain accuracy is also low. By contrast, we found in preliminary experiments that post-attack retain accuracy for many of the baselines remained high. We note that this is acceptable in our threat model as long as the pre-attack retain accuracy is still high, since the goal is to avoid recovering performance on the forget domain after tampering attacks at any cost. However, this does mean benign users trying to fine-tune the model must ensure their data is not contaminated by forget set data, lest their fine-tuned model have poor performance. This could make the method harder to use in practice. Thus, maintaining high post-attack retain accuracy would be a useful direction for future work to explore.

Tamper-resistance alone cannot fully mitigate the risks of malicious AI use. While it raises the initial costs for adversaries, it can eventually be circumvented. Once open-weight models are released, they cannot be "unreleased," leaving any compromised defenses permanently vulnerable. Therefore, tamper-resistance should be considered a supplement to the broader effort of of improving the offense-defense balance of AI systems. Addressing these limitations will improve the robustness of LLMs to tampering and better support open-weight model developers.

## B METHOD DETAILS

### B.1 INITIAL WEAPONIZATION KNOWLEDGE RESTRICTION SAFEGUARD

Prior to tamper-resistance training, we install a safeguard that achieves surgical knowledge restriction on the target hazardous domain. Let $h_\theta(\mathcal{D})$ denote the distribution of post-decoder layer residual stream activations for input sequences sampled from some data distribution $\mathcal{D}$ and model weights $\theta$. We define $\texttt{rand\_hashed}(x)$ for some input sequence $x$, which returns fixed Gaussian-sampled vectors that are chosen via hashing the corresponding input token for each residual stream index of $x$ in $\theta$. As a proxy for scrubbing target representations according to downstream task labels, we propose a weaponization knowledge restriction safeguard termed *Random Mapping*, which maps $h_\theta(\mathcal{D}_{\text{TR}})$ to random noise as follows:

$$\min_\theta \mathbb{E}_{x \sim \mathcal{D}_{\text{TR}}} \left[ 1 - \left| \frac{h_\theta(x) \cdot \texttt{rand\_hashed}(x)}{\|h_\theta(x)\|\|\texttt{rand\_hashed}(x)\|} \right| \right] + \mathcal{L}_{\text{LM}}(\theta; \mathcal{D}_{\text{retain}}) \qquad (3)$$

The objective of Equation 3 maximizes cosine similarity between row vectors in the residual stream in every layer of the LLM from $h(\mathcal{D}_{\text{TR}})$ and the hashed random vectors from $\texttt{rand\_hashed}(\cdot)$. By providing each token's residual stream a unique random vector to push toward, the loss encourages a "re-mapping" of token representations from $\mathcal{D}_{\text{TR}}$ to the noised vectors. We include an additional term for preserving performance on $\mathcal{D}_{\text{retain}}$ via the language-modelling cross-entropy loss $\mathcal{L}_{\text{LM}}$. We show the performance of the raw *Random Mapping* safeguard as an ablation in Table 5, listed as "Excl. Adv. Training."

### B.2 DESIGNING THE TAMPER-RESISTANCE LOSS

**Weaponization knowledge restriction.** For weaponization knowledge restriction, we summarize our intuition for ideal tamper-resistance loss design, corroborated by our empirical findings in Figure 3: *we seek to flatten the adversary's loss at a high value, rather than simply raise its y-intercept.*

We choose the tamper-resistance loss $\mathcal{L}_{\text{TR}}$ as an entropy loss to be maximized during the adversary's cross-entropy fine-tuning trajectory, since maximizing entropy would impede the adversary's cross entropy loss from decreasing during fine-tuning. In other words, we wish to obtain $\theta$ such that *after* an adversary performs a fine-tuning attack on $\theta$ via a cross-entropy loss, entropy is still high. We

find that this formulation achieves the desired flattening behavior, and we depict the difference in flattening between the choosing $\mathcal{L}_{TR}$ to be a negative cross-entropy loss and negative entropy loss in Figure 3. In the lefthand plot, where $\mathcal{L}_{TR}$ is a cross-entropy loss, loss only increases in the first inner loop step. In the righthand plot, where $\mathcal{L}_{TR}$ is a negative entropy loss, entropy is eventually maximized in all inner loop steps. Figure 5 also demonstrates the generalization of the flat adversary loss behavior beyond the length of the simulated fine-tuning trajectories during TAR.

**Harmful request refusal.** For harmful request refusal, we choose $\mathcal{L}_{TR}$ to be the DPO loss (Rafailov et al., 2023), which works as follows. Given a DPO dataset containing pairs of rejected and refusal completions, the sampled `attack` performs SFT on rejected completions, and the tamper-resistance loss $\mathcal{L}_{TR}$ is a DPO loss computed on the pair chosen and rejected completions on parameter coordinates along the `attack` trajectory. This encourages TAR to find an initialization $\theta$ such that after a harmful fine-tuning attack, the model still prefers refusal completions over harmful completions when given an harmful prompt. While this does not necessarily encourage a flat adversary loss, we find empirically in Table 2 that this formulation increases the average `safety_metric`$(\theta_G)$ defined in Section 3.2 after fine-tuning attacks on harmful data, detailed in Section 5.

We also observe that the length of the simulated adversary SFT trajectory during training affects test-time generalization in both Figure 5 and Appendix C.2. In particular, larger values of $K$ result in increased tamper-resistance for longer SFT attacks. However, to maintain reasonable runtime efficiency, we need a more efficient sampling technique than simply running $K$ independent trajectories of varying length for every outer-loop step in Algorithm 1, which we describe in Section B.3.

### B.3 EFFICIENTLY SAMPLING FINE-TUNING ATTACKS

Optimizing Equation 1 with gradient descent requires simulating $K$ tampering attacks for each tamper-resistance optimizer update, which is prohibitively expensive to run when the sampled `attack` performs SFT and $\theta$ contains billions of parameters. Inspired by prior work on snapshot ensembles Huang et al. (2017), we leverage an efficiency trick: we can reuse the coordinates along steps of a single adversary fine-tuning trajectory of length $K$ to obtain $K - 1$ additional (though non-independent) trajectories of increasing length. Using this trick, we collect all $K$ parameter coordinates along the trajectory into a single batch for computing the tamper-resistance losses, effectively sampling `attack` from $\mathcal{A}_{train}$ non-IID. To further improve runtime efficiency, we do not compute the tamper-resistance loss $\mathcal{L}_{TR}$ on all $K$ steps and instead sub-sample coordinates along the trajectory for computing $\mathcal{L}_{TR}$ within an adversary batch, for example every 4 adversary optimization steps. Additionally, we reduce variance in the tamper-resistance gradient by computing the tamper-resistance loss at each inner loop step on the same held-out batch, denoted as $x_{TR}$ in Algorithm 1.

### B.4 IMPLEMENTATION DETAILS AND RESOURCE REQUIREMENTS

We perform TAR training on Llama-3-8B-Instruct (Llama Team, AI @ Meta, 2024) with 8 NVIDIA 80GB A100 GPUs, leveraging distributed training via FSDP (Ren et al., 2021; Rajbhandari et al., 2020; Zhao et al., 2023). We use ZeRO Stage 3 from DeepSpeed (Rajbhandari et al., 2020), which shards optimizer states, gradients, and parameters during training. While the efficiency trick in Appendix B.3 improves runtime, we note additional considerations for conserving GPU memory.

First, simulating fine-tuning attacks that require additional state (e.g., momentum) in the inner loop of Algorithm 1 requires initializing a fresh optimizer for every outer loop iteration. Since we use an outer-loop optimizer that also requires maintaining state (ScheduleFree AdamW (Defazio et al., 2024)), we move the outer loop optimizer to the CPU before instantiating inner-loop optimizers.

Second, first-order meta-learning in smaller models can typically be implemented by running multiple forward passes for each inner loop iteration, averaging losses, then backpropagating on the averaged loss term. However, because each inner-loop tamper-resistance loss term ($\mathcal{L}_{TR}$ in Algorithm 1) is computed on a separate forward pass, this requires maintaining $K$ computation graphs in memory. Since this is infeasible on reasonable hardware for LLMs with billions of parameters, we circumvent this inefficiency by accumulating tamper-resistance gradients in a separate data structure ($g_{TR}$ in Algorithm 1). We note that this can be done without using additional `all-gather` and `reduce-scatter`

distributed operations, since tamper-resistance gradient accumulation and application to the pre-inner loop model parameters ($\theta_{i-1}$ in Algorithm 1) can be computed solely on sharded gradients.

## C  ADDITIONAL EXPERIMENTS

### C.1  BENIGN FINE-TUNING

|  |  | WMDP Forget ($\downarrow$) | Benign Domain ($\uparrow$) |
|---|---|---|---|
| TAR-Bio | Pre-SFT | 24.0 | 56.7 |
|  | Post-SFT | 24.0 | 64.2 |
| TAR-Cyber | Pre-SFT | 26.3 | 56.8 |
|  | Post-SFT | 27.5 | 62.2 |

Table 3: Average accuracy on MMLU economics subjects and WMDP Forget subjects for our Llama-3-8B TAR models safeguarded against hazardous biosecurity and cybersecurity knowledge, before and after fine-tuning on benign economics data (Li et al., 2024). Results indicate that models safeguarded with TAR still preserve benign fine-tunability.

An important property of open-weight models is that they can be fine-tuned to improve performance on custom data or in specific domains. Thus, ideal tamper-resistant safeguards should allow continued fine-tuning of a model while preserving the safeguard. We evaluate whether TAR models can be fine-tuned on data unrelated to the safeguard using economics as an example domain. Using TAR models with biosecurity and cybersecurity safeguards, we perform supervised fine-tuning on the WMDP auxiliary economics corpora (Li et al., 2024). We fine-tune models for $800$ steps using a learning rate of $2 \times 10^{-6}$ and a batch size of $64$. For the test set, we evaluate average accuracy on the corresponding MMLU High School Macroeconomics and Microeconomics subjects before and after fine-tuning. To confirm that the safeguard remains tamper-resistant in this setting, we also evaluate accuracy on corresponding WMDP subjects before and after fine-tuning.

We show the results of this evaluation in Table 3. We find that accuracy on economics questions can be improved by 7.5 percentage points without recovering significant hazardous knowledge. This illustrates that strong tamper-resistance can be compatible with benign model editing. Additional experiment details are discussed in Appendix D.

### C.2  VARYING THE TRAIN-TIME INNER-LOOP LENGTH $K$

Recall that via the efficiency trick discussed in Appendix B.3, a single inner loop trajectory of length $K$ during TAR returns the $K$ sampled attacks in Algorithm 1. We compare the test-time loss robustness as we vary the length of the inner loop $K$ during TAR, running fine-tuning attacks for $1{,}000$ steps on a held-out forget dataset for biosecurity weaponization (Adversary 8 in Table 9). For each value of $K$, we observe a plateau in the test loss that drops off at later steps as $K$ increases. This suggests that the robustness of TAR improves as the inner loop length increases. Prior work also corroborates that increasing the inner-loop length during meta-learning increases test-time generalization (Henderson et al., 2023). We note the contrast to conventional meta-learning methods mentioned in Section 4, in which typical meta-learning applications seek optimality after as few test-time steps as possible (Nichol et al., 2018; Finn et al.,
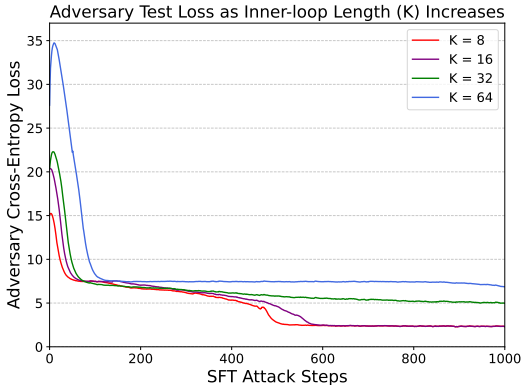


Figure 6: Comparison of a $1{,}000$-step SFT attack against TAR with the inner-loop length $K$ varied between $\{8, 16, 32, 64\}$. Test-time loss plateau magnitude and duration tend to increase as $K$ increases.

16

2017). Here, our results suggest that the TAR
objective is incentivized to run with as many
inner loop steps as possible. We find that $K = 64$ provides significant robustness to the range of
adversaries we consider in Section 5 and Appendix E.1, while balancing computational efficiency as
discussed in Appendix B.3.

## C.3 TARGETED PATCHING OF VULNERABILITIES

| Adversary Setup | Post-Attacks | |
| --- | --- | --- |
| | All Advs. ($\downarrow$) | R$\rightarrow$F ($\downarrow$) |
| Excl. R$\rightarrow$F Adversary | 50.2 | 63.9 |
| Incl. R$\rightarrow$F Adversary | **31.3** | **24.7** |

Table 4: Pre-Attack and Post-Attack scores when excluding and including the "Retain-set SFT
followed by Biosecurity SFT adversary" (R$\rightarrow$F) described in Appendix D.2. Average tamper-
resistance across all 28 adversaries improves by nearly 20 percentage points when including the R$\rightarrow$F
adversary in TAR, and improves by nearly 40 percentage points for the R$\rightarrow$F adversary at an LR of
$2 \times 10^{-5}$.

In Section D.2, we discuss the set of adversaries sampled during TAR for each weaponization
knowledge domain. However, we included the "Retain-set SFT followed by Weaponization-domain
SFT" (R$\rightarrow$F in Table 4) adversary to improve robustness to specific adversaries that broke the defense
in earlier iterations of the method that did not sample this adversary. We find in Table 4 that including
the R$\rightarrow$F adversary significantly improves robustness to more test-time attacks, indicated by an
improvement in Post-Attack Forget error by over 20.0 percentage points. Importantly, we find that
including the 64-step R$\rightarrow$F adversary at train-time improves test-time robustness at for the 1000-step
R$\rightarrow$F adversary by nearly 40.0 percentage points, indicating that target patching of vulnerabilities is
possible within TAR.

## C.4 ABLATIONS

| Ablation | Pre-Attack | | Post-Attacks (Avg) |
| --- | --- | --- | --- |
| | Retain ($\uparrow$) | Forget ($\downarrow$) | Forget ($\downarrow$) |
| No Defense | 67.3 | 70.5 | 70.5 |
| Excl. $\ell_2$ term from $\mathcal{L}_{\text{retain}}$ | 49.6 | 26.6 | 38.7 |
| Excl. Adv. Training | 59.7 | 27.3 | 59.4 |
| Excl. Initial Safeguard | 62.5 | 47.3 | 46.9 |
| TAR | 54.9 | **24.0** | **31.3** |

Table 5: Ablations for primary components of TAR: (1) the $\ell_2$ retain loss term, (2) the adversarial
training phase, (3) the initial model safeguard phase. We find that all 3 components are critical for
the high tamper-resistance that TAR achieves, indicated by an average Post-Attack Forget accuracy
decrease across all ablations compared to TAR of 17.0%.

**Retain loss $\ell_2$ representation engineering term.** We observe in Table 5 that retain-set MMLU
performance is worse when excluding representation engineering (Zou et al., 2023a) retain loss from
TAR. When including the term, retain MMLU improves by at least 5.0 percentage points.

**Including the initial safeguard.** In Table 5, we examine the impact of incorporating the *Random
Mapping* safeguard step prior to the adversarial training phase during TAR. The *Random Mapping*
safeguard in isolation achieves a near-random chance Pre-Attack Forget accuracy of 27.3. However,
it is susceptible to fine-tuning attacks similar to other baselines in Table 1, indicated by a higher
Post-Attack Forget accuracy of 59.4. When including the tamper-resistance adversarial training phase,

we observe significantly increased tamper-resistance as the Post-Attack Forget accuracy decreases by nearly 28 percentage points.

We also examine the impact of excluding the initial safeguarding step ("Excl. Initial Safeguard" in Table 5), finding that Pre-Attack and Post-Attack Forget accuracies are both higher without the initial safeguard. While the tamper-resistance adversarial training phase achieves moderate tamper-resistance without the initial safeguard step, we find that including the *Random Mapping* phase improves downstream tamper-resistance by 15.6 percentage points.

| $\mathcal{L}_{\text{TR}}$ Weighting | Pre-Attack | | Post-Attacks (Avg) |
|---|---|---|---|
| | Retain ($\uparrow$) | Forget ($\downarrow$) | Forget ($\downarrow$) |
| $\lambda_{\text{TR}} = 1.0$ | 62.5 | 29.3 | 40.8 |
| $\lambda_{\text{TR}} = 4.0$ | 54.9 | **24.0** | **31.3** |

Table 6: Pre-Attack and Post-Attack scores when varying the tamper-resistance loss weighting, $\lambda_{\text{TR}}$. Tamper-resistance improves by nearly $10.0\%$ when increasing $\lambda_{\text{TR}}$ from $1.0$ to $4.0$. The retain loss weight $\lambda_{\text{retain}}$ is fixed at $1.0$ for both settings.

**Varying the tamper-resistance loss scale $\lambda_{\text{TR}}$.** We compare the downstream robustness of TAR when varying the tamper-resistance loss weighting $\lambda_{\text{TR}}$ between $1.0$ and $4.0$ in Table 6. We observe that when setting $\lambda_{\text{TR}} = 1.0$, TAR maintains high retain MMLU accuracy at $62.5$ percentage points, with moderate tamper-resistance indicated by a Post-Attack Forget accuracy of $40.8$. Further increasing $\lambda_{\text{TR}}$ to $4.0$ in our final TAR model results in a significantly improved Post-Attack Forget Accuracy of $31.3$, with a partial decrease in Retain MMLU to $54.9$. When varying $\lambda_{\text{TR}}$, we keep $\lambda_{\text{retain}}$ constant; thus, our results indicate a clear way to increase downstream tamper-resistance by increasing the weighting of the tamper-resistance gradient during TAR, reflecting a balance between tamper-resistance and capabilities similar to the robustness-performance tradeoff for adverarial robustness in vision models.

### C.5 DPO Tamper-Resistance Loss Win-Rate during TAR

In Figure 7, we plot the DPO win-rate during harmful SFT attack trajectories during the adversarial training phase of TAR. We find that the outer-loop DPO loss steadily decreases, which corresponds to the average inner-loop win-rate of refusal completions over rejected completions steadily increasing over the 100 outer loop steps. Our results demonstrate that TAR is able to satisfy complex tamper-resistance losses after fine-tuning. We believe that this is a useful feature of the method, enabling TAR to adapt to other potentially useful objective functions that correspond to downstream robustness.



Figure 7: The development of inner-loop DPO win-rates during harmful SFT attack inner loops (red), over the course of tamper-resistance training for TAR. The outer loop win-rate (blue) depicts the average win-rate across inner loops over the course of tamper-resistance training. We observe that by the end of training, the win-rate for refusal completions becomes completely flat near the optimal win-rate value of $1.0$. TAR is run for $N = 100$ outer loop steps, and the average DPO Win-Rate on the $y$-axis is bounded between $0$ and $1$.

## D Experiment Details

### D.1 Weaponization Domain Proxy Dataset Details

**Biosecurity.** We use a synthetically labeled partition of the Pile (Gao et al., 2020) that filters for relevance to biology and the Camel AI Biology dataset (Li et al., 2023). We generate synthetic labels for Pile token sequences using openchat-3.5 (Wang et al., 2023), categorizing

them as belonging to "Cellular Biology" or not.

This process yields 49,984 samples: 7,558 for the Forget-set (Pile-bio Forget) and 42,426 for the Retain-set (Pile-bio Retain). Concurrently, we pack entries from the Camel AI Biology dataset to the truncation-enabled 256 tokenization limit, resulting in 54,258 samples of about 188 words each (Camel-bio Forget). We apply the same procedure to our held-out hazardous biology dataset (identical to the WMDP biosecurity Forget-set), producing 598,933 samples of similar length (OOD Forget).

**Chemical Security.** We use a private forget dataset containing text sequences about hazardous chemical security content (Chem Forget).

**Cybersecurity.** We scrape CTF writeups on CTFtime (CTFtime, 2024) that are numbered between 1 and 39181, collecting cybersecurity writeups written as recently as 2024. We filter to keep writeups that contain more than 150 characters. As a result of filtering and HTTP errors while scraping, our resulting forget dataset contains slightly over 18k samples (Cyber Forget).

### D.2 Train-time Settings and Adversaries

**Weaponization knowledge restriction.** For each weaponization knowledge restriction domain, we have a corresponding retain dataset $\mathcal{D}_{\text{retain}}$, comprised of a mix of data from the Pile-bio Retain set and Magpie-Align instruction-tuning dataset (Xu et al., 2024), which we refer to as the Retain-set. Referencing the datasets described in Appendix D.1, we specify the adversary tamper-resistance datasets that constitute the data used by attacks sampled from $\mathcal{A}_{\text{train}}$ during TAR training, as well as the attack setup for each domain as follows:

- **Biosecurity:** We simulate six adversaries from the following three setups: Pile-Bio Forget-set SFT, Camel-Bio Forget-set SFT, and Retain-set SFT followed by Pile-Bio Forget-set SFT (R→F), where the switching point between Retain-set and Pile-Bio Forget-set SFT within the 64-step length trajectory is sampled from a beta distribution $\beta(6.0, 3.0)$. LRs are sampled from $\{2 \times 10^{-5}, 4 \times 10^{-5}\}$.
- **Chemical Security:** We simulate nine adversaries from the following three setups: Chem Forget-set SFT, Retain-set SFT, and Retain-set followed by Chem Forget-set SFT, using the same switching-point sampling scheme as in the Biosecurity setting. LRs are sampled from $\{2 \times 10^{-6}, 2 \times 10^{-5}, 4 \times 10^{-5}\}$.
- **Cybersecurity:** We simulate four adversaries from the following two setups: Cyber Forget-set SFT, and Retain-set SFT followed by Cyber Forget-set SFT, using the same switching point sampling scheme as in the Chemical Security and Biosecurity settings. LRs are sampled from $\{2 \times 10^{-5}, 4 \times 10^{-5}\}$.

For each weaponization knowledge domain, we create 80-20 splits for adversary and held-out data of the corresponding forget sets, respectively. For Biosecurity, which uses multiple forget datasets, this involves creating 80-20 splits for each dataset, then combining the corresponding splits. The adversary data splits are used for sampled attacks from $\mathcal{A}_{\text{train}}$, whereas the held-out split is used for computing tamper-resistance losses. The held-out splits for each domain correspond to $\mathcal{D}_{\text{TR}}$ in Section 4. We use minibatches from a held-out dataset for computing tamper-resistance losses rather than cycling through a single dataset, following the recommendation of Nichol et al. (2018), in which first-order meta-learning without properly held-out minibatches caused a performance degradation.

All train-time adversary setups are tabulated in Table 7, where F-Pile, F-Chem, F-Cyber denote the respective datasets described in Appendix D.1, and Retain denotes the mixed Pile-bio and Magpie-Align Retain-set described in Appendix D.2. We use R→F to label the adversaries that perform Retain-set SFT followed by Forget-set SFT. The final column is an abbreviation for Finetuning Paradigm and indicates whether the SFT setup used full parameter finetuning or parameter-efficient fine-tuning (PEFT) via LoRA adapters (Hu et al., 2021).

Table 7: Train-time adversary setups for weaponization knowledge restriction of Biosecurity, Chemical Security, and Cybersecurity.

| Adversary | Dataset | Opt. Steps ($K$) | Optimizer | LR | LR Schedule | Batch Size | FT Paradigm |
|-----------|---------|------------------|-----------|-----|-------------|------------|-------------|
| **Biosecurity Weaponization Restriction** | | | | | | | |
| Adv 1 | F-Pile | 64 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 2 | F-Pile | 64 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 3 | F-Camel | 64 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 4 | F-Camel | 64 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 5 | R→F | 64 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 6 | R→F | 64 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| **Chemical Security Weaponization Restriction** | | | | | | | |
| Adv 1 | F-Chem | 64 | AdamW | $2 \times 10^{-6}$ | No Warmup | 64 | Full Parameter |
| Adv 2 | F-Chem | 64 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 3 | F-Chem | 64 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 4 | Retain | 64 | AdamW | $2 \times 10^{-6}$ | No Warmup | 64 | Full Parameter |
| Adv 5 | Retain | 64 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 6 | Retain | 64 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 7 | R→F | 64 | AdamW | $2 \times 10^{-6}$ | No Warmup | 64 | Full Parameter |
| Adv 8 | R→F | 64 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 9 | R→F | 64 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| **Cybersecurity Weaponization Restriction** | | | | | | | |
| Adv 1 | F-Cyber | 64 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 2 | F-Cyber | 64 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 3 | R→F | 64 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 4 | R→F | 64 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |

**Harmful request refusal.** For harmful request refusal, we choose the retain dataset $\mathcal{D}_{\text{retain}}$ to be the Magpie-Align instruction-tuning dataset (Xu et al., 2024). We sample train-time adversaries that perform $K = 64$ steps of SFT on rejected completions of the Anthropic-HH-RLHF dataset (Bai et al., 2022) and vary the learning rate within $\{2 \times 10^{-6}, 2 \times 10^{-5}, 4 \times 10^{-5}\}$. We depict this list of adversaries in Table 8.

While the rejected completions from Anthropic-HH-RLHF constitute the data used for sampled attacks from $\mathcal{A}_{\text{train}}$ for harmful request refusal, we compute the tamper-resistance loss $\mathcal{L}_{\text{TR}}$ as follows. Since $\mathcal{L}_{\text{TR}}$ is the DPO loss (Rafailov et al., 2023) in this setting, we use the base model weights $\theta$ as the reference model and sample harmful and benign completions from a modified test split of Anthropic-HH-RLHF (Bai et al., 2022), where rejected completions are replaced with refusals (Cai et al., 2024). To avoid keeping the base model weights in memory and speed up training, we precompute the reference model DPO log-probabilities for the full Anthropic-HH-RLHF dataset before training. To summarize, we have that the sampled train-time adversaries perform SFT on the *rejected* completions from Anthropic-HH-RLHF, and tamper-resistance DPO loss is computed on the corresponding modified refusal completions; the modified refusal completions in this setting correspond to $\mathcal{D}_{\text{TR}}$ in Section 4.

In practice, we perform an additional 100 steps of supervised fine-tuning on the Magpie-Align dataset to improve the benign capabilities performance of the TAR refusal model in Table 2.

Table 8: Train-time adversary red-teaming setups harmful request refusal. "A-HH-Rejected" in the Dataset column corresponds to the adversary dataset comprised of rejected completions from the Anthropic-HH-RLHF dataset.

| Adversary | Dataset | Opt. Steps ($K$) | Optimizer | LR | LR Schedule | Batch Size | FT Paradigm |
|---|---|---|---|---|---|---|---|
| **Harmful Request Refusal** | | | | | | | |
| Adv 1 | A-HH-Rejected | 64 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 2 | A-HH-Rejected | 64 | AdamW | $2 \times 10^{-6}$ | No Warmup | 64 | Full Parameter |
| Adv 3 | A-HH-Rejected | 64 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |

# E RED TEAMING DETAILS

## E.1 WEAPONIZATION KNOWLEDGE RESTRICTION

In Table 9, we list all *test-time* adversary setups for recovering Biosecurity, Chemical Security, and Cybersecurity Weaponization knowledge.

For Biosecurity, we examine post-attack forget accuracy after fine-tuning for $1,000$ steps on three data distributions: the Pile-bio Forget set and the Retain-set used during *Random Mapping* and TAR, and an OOD-Forget set mentioned in Appendix D.1, which is unseen during tamper-resistance training. We conduct an additional attack using $5,000$ optimization steps on this unseen forget distribution to compose our strongest adversary. We leverage a consistent setup for Chemical Security and Cybersecurity, substituting the Pile-bio Forget set for the respective Chem and Cyber Forget sets. We exclude the unseen forget distribution adversaries for a total of $24$ attacks, as we were unable to obtain an equally high quality held-out dataset for Chemical Security and Cybersecurity.

For all subjects, we compare the effect of different optimizers (AdamW, Adadelta, and Stochastic Gradient Descent with Nesterov Momentum, Schedule Free AdamW) (Zeiler, 2012; Kingma & Ba, 2017; Xie et al., 2023; Defazio et al., 2024), learning rates ($2 \times 10^{-6}, 2 \times 10^{-5}, 4 \times 10^{-5}$), and learning rate scheduling techniques (SGDR and 30 steps of linear warmup) (Loshchilov & Hutter, 2016).

Similar to Table 7, we use F-Pile, F-Chem, F-Cyber, and OOD-F in the Dataset Column to denote the respective datasets described in Appendices E.1. At test-time, we use the Pile-bio Retain set as the global Retain-set adversary, labeled as Retain. The R→F adversary at test-time also differs from the train-time version: we perform Forget-set SFT for 40% of the optimization steps, followed by Retain-set SFT for the remaining 60%. We found this combination to be a potent attack that broke intermediate versions of the method, as described in Appendix C.3.

In the Optimizer column, Schedule Free is an abbreviation of Schedule Free AdamW and SGD Nesterov is an abbreviation of SGD with Nesterov Momentum. In cases where the adversary used parameter-efficient fine-tuning (PEFT) via LoRA adapters (Hu et al., 2021), we used a LoRA config with an attention dimension, or rank, of 16, a LoRA alpha value of 32, a LoRA dropout of 0.05, on target linear modules:

```
{'up_proj', 'down_proj', 'gate_proj', 'q_proj', 'k_proj', 'v_proj', 'o_proj'}.
```

Table 9: Test-time adversary red-teaming setups for weaponization knowledge restriction of Biosecurity, Chemical Security, and Cybersecurity.

| Adversary | Dataset | Opt. Steps ($K$) | Optimizer | LR | LR Schedule | Batch Size | FT Paradigm |
|---|---|---|---|---|---|---|---|
| **Biosecurity Weaponization Restriction** | | | | | | | |
| Adv 1 | OOD-F | 5000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| | | | | | | | *Continued on next page* |

21

Under review as a conference paper at ICLR 2025

**Table 9 – continued from previous page**

| Adversary | Dataset | Opt. Steps ($K$) | Optimizer | LR | LR Schedule | Batch Size | FT Paradigm |
|---|---|---|---|---|---|---|---|
| Adv 2 | OOD-F | 5000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 3 | F-Pile | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 4 | F-Pile | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 5 | Retain | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 6 | Retain | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 7 | OOD-F | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 8 | OOD-F | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 9 | R→F | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 10 | R→F | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 11 | F-Pile | 1000 | Adadelta | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 12 | F-Pile | 1000 | Adadelta | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 13 | F-Pile | 1000 | Schedule Free | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 14 | F-Pile | 1000 | Schedule Free | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 15 | F-Pile | 1000 | SGD Nesterov | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 16 | F-Pile | 1000 | SGD Nesterov | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 17 | F-Pile | 1000 | AdamW | $2 \times 10^{-6}$ | No Warmup | 64 | Full Parameter |
| Adv 18 | F-Pile | 1000 | AdamW | $2 \times 10^{-6}$ | 30 Steps Warmup | 64 | Full Parameter |
| Adv 19 | F-Pile | 1000 | AdamW | $2 \times 10^{-5}$ | 30 Steps Warmup | 64 | Full Parameter |
| Adv 20 | F-Pile | 1000 | AdamW | $4 \times 10^{-5}$ | 30 Steps Warmup | 64 | Full Parameter |
| Adv 21 | F-Pile | 1000 | AdamW | $2 \times 10^{-5}$ | SGDR | 64 | Full Parameter |
| Adv 22 | F-Pile | 1000 | AdamW | $4 \times 10^{-5}$ | SGDR | 64 | Full Parameter |
| Adv 23 | F-Pile | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 32 | Full Parameter |
| Adv 24 | F-Pile | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 32 | Full Parameter |
| Adv 25 | F-Pile | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 128 | Full Parameter |
| Adv 26 | F-Pile | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 128 | Full Parameter |
| Adv 27 | F-Pile | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | PEFT |
| Adv 28 | F-Pile | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | PEFT |
| **Chemical Security Weaponization** | | | | | | | |
| Adv 1 | F-Chem | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 2 | F-Chem | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 3 | Retain | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 4 | Retain | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 5 | R→F | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 6 | R→F | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 7 | F-Chem | 1000 | Adadelta | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 8 | F-Chem | 1000 | Adadelta | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 9 | F-Chem | 1000 | ScheduleFree | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 10 | F-Chem | 1000 | ScheduleFree | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 11 | F-Chem | 1000 | SGD Nesterov | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 12 | F-Chem | 1000 | SGD Nesterov | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 13 | F-Chem | 1000 | AdamW | $2 \times 10^{-6}$ | No Warmup | 64 | Full Parameter |
| Adv 14 | F-Chem | 1000 | AdamW | $2 \times 10^{-6}$ | 30 Steps Warmup | 64 | Full Parameter |
| Adv 15 | F-Chem | 1000 | AdamW | $2 \times 10^{-5}$ | 30 Steps Warmup | 64 | Full Parameter |
| Adv 16 | F-Chem | 1000 | AdamW | $4 \times 10^{-5}$ | 30 Steps Warmup | 64 | Full Parameter |
| Adv 17 | F-Chem | 1000 | AdamW | $2 \times 10^{-5}$ | SGDR | 64 | Full Parameter |
| Adv 18 | F-Chem | 1000 | AdamW | $4 \times 10^{-5}$ | SGDR | 64 | Full Parameter |
| Adv 19 | F-Chem | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 32 | Full Parameter |
| Adv 20 | F-Chem | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 32 | Full Parameter |
| Adv 21 | F-Chem | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 128 | Full Parameter |
| Adv 22 | F-Chem | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 128 | Full Parameter |
| Adv 23 | F-Chem | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | PEFT |
| Adv 24 | F-Chem | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | PEFT |
| **Cybersecurity Weaponization Restriction** | | | | | | | |
| Adv 1 | F-Cyber | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |

Continued on next page

**Table 9 – continued from previous page**

| Adversary | Dataset | Opt. Steps ($K$) | Optimizer | LR | LR Schedule | Batch Size | FT Paradigm |
|---|---|---|---|---|---|---|---|
| Adv 2 | F-Cyber | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 3 | Retain | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 4 | Retain | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 5 | R→F | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 6 | R→F | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 7 | F-Cyber | 1000 | Adadelta | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 8 | F-Cyber | 1000 | Adadelta | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 9 | F-Cyber | 1000 | ScheduleFree | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 10 | F-Cyber | 1000 | ScheduleFree | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 11 | F-Cyber | 1000 | SGD Nesterov | $2 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 12 | F-Cyber | 1000 | SGD Nesterov | $4 \times 10^{-5}$ | No Warmup | 64 | Full Parameter |
| Adv 13 | F-Cyber | 1000 | AdamW | $2 \times 10^{-6}$ | No Warmup | 64 | Full Parameter |
| Adv 14 | F-Cyber | 1000 | AdamW | $2 \times 10^{-6}$ | 30 Steps Warmup | 64 | Full Parameter |
| Adv 15 | F-Cyber | 1000 | AdamW | $2 \times 10^{-5}$ | 30 Steps Warmup | 64 | Full Parameter |
| Adv 16 | F-Cyber | 1000 | AdamW | $4 \times 10^{-5}$ | 30 Steps Warmup | 64 | Full Parameter |
| Adv 17 | F-Cyber | 1000 | AdamW | $2 \times 10^{-5}$ | SGDR | 64 | Full Parameter |
| Adv 18 | F-Cyber | 1000 | AdamW | $4 \times 10^{-5}$ | SGDR | 64 | Full Parameter |
| Adv 19 | F-Cyber | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 32 | Full Parameter |
| Adv 20 | F-Cyber | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 32 | Full Parameter |
| Adv 21 | F-Cyber | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 128 | Full Parameter |
| Adv 22 | F-Cyber | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 128 | Full Parameter |
| Adv 23 | F-Cyber | 1000 | AdamW | $2 \times 10^{-5}$ | No Warmup | 64 | PEFT |
| Adv 24 | F-Cyber | 1000 | AdamW | $4 \times 10^{-5}$ | No Warmup | 64 | PEFT |

### E.2 HARMFUL REQUEST REFUSAL

For the harmful request refusal setting, we conduct $5$ test-time adversary attacks that perform SFT for $10$ epochs on a held-out toxicity dataset called Toxic-DPO v0.2, on each of the settings in Table 10. The dataset contains $541$ user-assistant chat interactions where the assistant complies with harmful instructions.

Table 10: Test-time adversary red-teaming setups for harmful request refusal. "ToxicDPO" in the Dataset column refers to the ToxicDPOv0.2 dataset containing harmful chat completions.

| Adversary | Dataset | Epochs | Optimizer | LR | LR Schedule | Batch Size | FT Paradigm |
|---|---|---|---|---|---|---|---|
| **Harmful Request Refusal** | | | | | | | |
| Adv 1 | ToxicDPO | 10 | AdamW | $1 \times 10^{-5}$ | 10 Steps Warmup | 32 | Full Parameter |
| Adv 2 | ToxicDPO | 10 | AdamW | $1 \times 10^{-5}$ | No Warmup | 32 | Full Parameter |
| Adv 3 | ToxicDPO | 10 | AdamW | $1 \times 10^{-5}$ | No Warmup | 16 | Full Parameter |
| Adv 4 | ToxicDPO | 10 | AdamW | $2 \times 10^{-5}$ | No Warmup | 32 | Full Parameter |
| Adv 5 | ToxicDPO | 10 | AdamW | $4 \times 10^{-5}$ | No Warmup | 32 | Full Parameter |

### E.3 ADDITIONAL WEAPONIZATION KNOWLEDGE RESTRICTION RED-TEAMING RESULTS

In Figure 8, we characterize the post-tampering resistance of TAR and the baseline safeguards with respect to the Chemical Security and Cybersecurity domains. Similar to the Biosecurity domain, baseline safeguards break under the majority of fine-tuning attacks, while TAR resists attacks from 23 out of 24 Chemical Security SFT adversaries and 21 out of 24 Cybersecurity SFT adversaries.
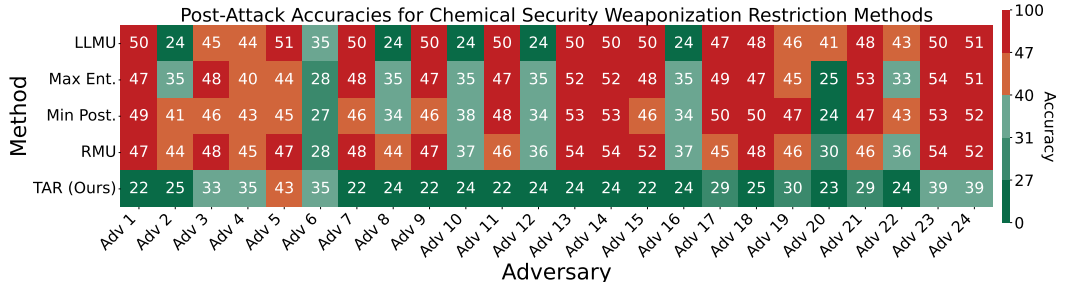
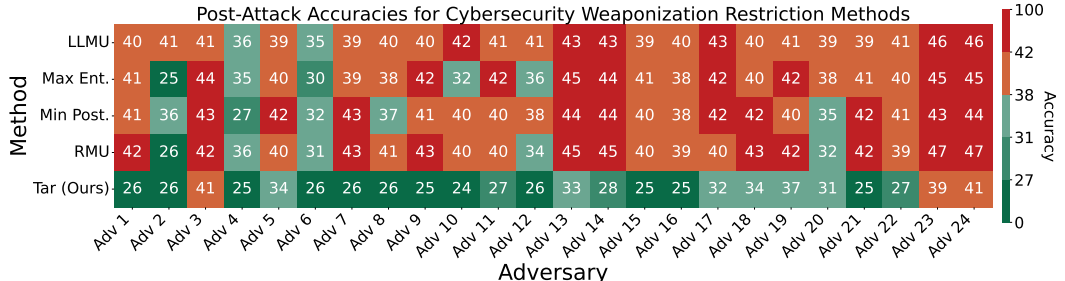| Model | Pre-Attacks | | Post-Attacks (Avg) |
|---|---|---|---|
| | MT-Bench ($\uparrow$) | ASR ($\downarrow$) | ASR ($\downarrow$) |
| Refusal Trained | 8.1 | 14.7 | 72.5 |
| R2D2 | 6.0 | 25.0 | 78.3 |
| RepNoise | 6.2 | 18.8 | 74.5 |
| RR | 8.0 | **1.4** | 84.8 |
| TAR (Ours) | 6.3 | 31.4 | **63.9** |

Table 11: Average Post-Attack HarmBench ASR, reported for TAR, Representation Rerouting (RR), and the Refusal Trained Llama-3-8B-Instruct model across 5 fine-tuning attacks depicted in Appendix E.2, as well as Pre-Attack MT-Bench and HarmBench ASR. TAR is more robust than other methods after tampering, while maintaining comparable MT-Bench performance. Note that Pre-Attack ASR is not a priority for us, as we focus on reducing ASR after tampering attacks. To improve both metrics, future work could consider combining tamper-resistance training with a strong baseline safeguard like RR. ASR values are percentages.



Figure 8: Red-teaming results for weaponization knowledge restriction in the Chemical Security and Cybersecurity domains. All numbers are percentages, with ideal defenses obtaining $25\%$. Red colors indicate that an attack recovered performance near the level of the No Defense baseline. We evaluate each defense against a diverse range of strong adversaries, described in Appendix E.1. Compared to prior safeguards, TAR greatly increases tamper-resistance for nearly all adversaries.

# F  BASELINE DETAILS

## F.1  WEAPONIZATION KNOWLEDGE RESTRICTION

**Max Entropy.**  Let $K$ be the set of all token-wise output probability distributions returned by a model $\theta$, where $k \in K$ corresponds to every position in the sequence. We maximize the average entropy of these discrete distributions in $K$ as follows:

$$\mathcal{L}_{\text{Max Entropy}} = \sum_{k \in \mathcal{K}} p_k \log(p_k)$$

24

This is equivalent to minimizing the average Kullback-Leibler (KL) divergence (Kullback & Leibler, 1951) between each $k$ and the discrete uniform distribution $u(x)$ over the vocabulary $V$. For Llama-3-8B-Instruct, $|V| = 128256$. Thus, this objective is upper-bounded by:

$$h(x) = -\log(u(x)) = \log(|V|) \approx 11.76$$

where $h(x)$ measures the Shannon information or self-information and $\log(x)$ has base $e$. We apply this objective for all elements in the Forget-set and perform standard cross-entropy on the Retain-set.

**Min Posterior.** The goal of the Min Posterior objective is to assign lower probabilities to true forget-set labels, essentially minimizing $-\log(1 - P(\text{label}))$. Let $p_i$ be the probability assigned to the true label for token $i$ and $\mathcal{V}$ be the model's vocabulary distribution. We define the Min Posterior objective function as follows:

$$\mathcal{L}_{\text{Min Posterior}} = -\frac{1}{|\mathcal{V}|} \sum_{i \in \mathcal{V}} \log(1 - p_i + \epsilon) \cdot \mathbb{I}[\log(p_i) \geq \tau]$$

where $\tau$ is the threshold for masking out target label logits (which we set to the negative maximum entropy of the vocabulary distribution, $-\log|\mathcal{V}|$) and $\mathbb{I}[\cdot]$ is the corresponding indicator function (1 if the condition is true, 0 otherwise). We include an optional $\epsilon = 1 \times 10^{-12}$ to help with numerical stability. Similar to the Max Entropy objective, we apply this objective for all elements in the Forget-set and perform standard cross-entropy on the Retain-set.

**RMU.** We adapt RMU's implementation from Li et al. (2024) with a learning rate of $5 \times 10^{-5}$ and 250 unlearning steps. We use the released WMDP's unlearning datasets for Biosecurity (Bio) and Cybersecurity (Cyber) unlearning, and our private hazardous chemistry dataset for Chemical Security (Chem) unlearning. We use unlearning coefficients of 20, 30, and 50 for Bio, Cyber, and Chem respectively. We use a retain coefficient of 700 on Wikitext Merity et al. (2016).

**LLMU.** We use a modified version of LLMU from Yao et al. (2023). Instead of computing the KL divergence to regularize retain-set logits towards the base frozen model, we employ a standard cross-entropy loss. This modification allows for memory-efficient execution on our hardware while maintaining comparable performance.

**Hyperparameter tuning.** Besides RMU, all baseline hyperparameters were chosen after a grid search across learning rates $\{3 \times 10^{-6}, 5 \times 10^{-6}, 8 \times 10^{-6}, 1 \times 10^{-5}\}$, optimization step count $\{600, 1000\}$, and warmup steps $\{0, 100\}$. We found that 600 optimization steps using the AdamW Schedule Free optimizer, at a learning rate of $1 \times 10^{-5}$, with 100 steps of linear warmup, and an effective batch size of 64 produced the best performance. For the Max Entropy, Min Posterior, and LLMU baselines, we train on the three corresponding forget datasets discussed in D.1. For these baselines, we modify our Biosecurity forget corpus to be a mixture of the Pile-bio and Camel-bio Forget corpora. We use the Pile-bio Retain-set as a global Retain-set for baseline training.

### F.2 HARMFUL REQUEST REFUSAL

**Representation Rerouting.** We use the Llama-3-8B-Instruct RR model from Zou et al. (2024), which uses a cosine distance loss to push representations for harmful inputs to become orthogonal to those of the base Llama-3-8B-Instruct model.

**R2D2.** We use the R2D2 model run on Zephyr-7B directly from Mazeika et al. (2024), which performs adversarial training against GCG attacks to increase jailbreak robustness.

**RepNoise.** We the RepNoise model run on Llama-2-7B directly from Rosati et al. (2024b), which uses a distributional loss to push representations for harmful inputs toward Gaussian noise.

### F.3 ADDITIONAL BASELINE COMPARISONS

**MLAC-AR.** Meta-Learned Adversarial Censoring (MLAC) (Henderson et al., 2023) was originally proposed to prevent BERT-style models from learning binary classification for gender bias data.

| Domain | Model | Pre-Attacks | | Post-Attacks (Avg) |
|---|---|---|---|---|
| | | Retain ($\uparrow$) | Forget ($\downarrow$) | Forget ($\downarrow$) |
| Biosecurity | Random | 25.0 | 25.0 | 25.0 |
| | MLAC-AR | 49.1 | 31.2 | 59.6 |
| | SOPHON-AR | 27.2 | 24.0 | 32.1 |
| | TAR (Ours) | **54.9** | 24.0 | 31.3 |
| Chemical Security | Random | 25.0 | 25.0 | 25.0 |
| | MLAC-AR | 47.8 | 29.9 | 32.3 |
| | SOPHON-AR | 23.3 | 26.2 | 28.1 |
| | TAR (Ours) | **57.9** | 29.2 | 27.6 |
| Cybersecurity | Random | 25.0 | 25.0 | 25.0 |
| | MLAC-AR | 36.0 | 26.6 | 30.9 |
| | SOPHON-AR | 24.4 | 24.6 | 25.4 |
| | TAR (Ours) | **55.8** | 26.3 | 29.5 |

Table 12: Additional baselines for MLAC-AR, an extension of the method in Henderson et al. (2023) to autoregressive LLMs, as well as SOPHON-AR from Deng et al. (2024), respectively. Despite extensive tuning, SOPHON-AR does not yield a usable model. Additionally, MLAC-AR has varying robustness and worse Retain MMLU performance.

Since the approach is not immediately applicable to LLMs, we extend MLAC in a variant we call autoregressive MLAC (MLAC-AR). Since MLAC in its original formulation calls for "task-blocking" via negating the adversary's loss during the inner loop of meta-learning, we implement this by negating the cross-entropy loss of an LLM fine-tuning adversary. However, we found that this approach diverges in performance across a variety of hyperparameters, and opted to further improve performance of the MLAC-AR baseline by clamping the maximum cross-entropy loss at the value of the maximum entropy of the output vocabulary distribution, $\log(\texttt{vocab\_size})$. We show results in Table 12, finding that MLAC-AR does not maintain sufficient benign capabilities performance nor uniform tamper-resistance across weaponization domains.

**SOPHON-AR.** In concurrent work, SOPHON (Deng et al., 2024) was introduced to prevent small diffusion models and image classifiers from learning specific data distributions. Similarly to MLAC-AR, we extend SOPHON to LLMs via SOPHON-AR, using the alternating retain loss and fine-tuning suppression loss formulation that the authors propose. Furthermore, we adapt the inverse cross-entropy loss from Deng et al. (2024) , which aims to boost convergence of the fine-tuning suppression process. We find in practice that despite heavy tuning, SOPHON-AR does not converge well enough to yield a usable Pre-Attack model in Table 12.