Improving Large Language Model Safety with Contrastive Representation Learning

Samuel Simko

ETH Zurich ssimko@student.ethz.ch

Bernhard Schölkopf

MPI for Intelligent Systems bs@tue.mpg.de

Abstract

Large Language Models (LLMs) are powerful tools with profound societal impacts, yet their ability to generate responses to diverse and uncontrolled inputs leaves them vulnerable to adversarial attacks. While existing defenses often struggle to generalize across varying attack types, recent advancements in representation engineering offer promising alternatives. In this work, we propose a defense framework that formulates model defense as a contrastive representation learning (CRL) problem. Our method finetunes a model using a triplet-based loss combined with adversarial hard negative mining to encourage separation between benign and harmful representations. Our experimental results across multiple models demonstrate that our approach outperforms prior representation engineering-based defenses, improving robustness against both input-space and embeddingspace attacks without compromising standard performance. 1

1 Introduction

In recent years, Large Language Models (LLMs) have proven to be powerful tools for general-purpose language understanding and generation (Minaee et al., 2024). They have had significant impact on software engineering (Hou et al., 2024), medicine (Zhou et al., 2024) and natural sciences (Nejjar et al., 2024).

However, their ability to respond to uncontrolled inputs comes with serious security risks (Geiping et al., 2024), as they can generate inappropriate, toxic, or harmful text (Deshpande et al., 2023). To mitigate this risk, various mechanisms have been developed to keep LLMs aligned with human values (Stiennon et al., 2022; Bai et al., 2022; Rafailov et al., 2024). Nonetheless, the most commonly used LLM systems often fail to protect against

¹Our code is available at https://github.com/samuelsimko/crl-llm-defense

Mrinmaya Sachan ETH Zurich

msachan@ethz.ch

Zhijing Jin

MPI & University of Toronto zjin@cs.toronto.edu

such behaviors (Zou et al., 2023; Chao et al., 2023). Developing safe and reliable defenses is therefore essential for minimizing societal risks associated with broad AI deployment.

The current state-of-the-art defenses against LLM jailbreaks focus either on creating effective pipelines around the model (Zeng et al., 2024b; Yuan et al., 2024b), or on finetuning the model weights (Yousefpour et al., 2025; Zhang et al., 2024). Among the latter, circuit breakers (Zou et al., 2024) are a promising approach, designed to disrupt the model's harmful inner representations, causing it to generate incoherent or nonsensical outputs rather than harmful content even under attacks.

In this work, motivated by the success of circuit breaking, we introduce a new approach to representation space safety engineering by building on concepts from contrastive learning (Khosla et al., 2021). Our triplet loss formulation represents a natural extension of circuit breaking, and allows for clearer separation of harmful and benign representations. Additionally, we adopt an improved training strategy based on adversarial hard negative mining (Casper et al., 2024). Overall, we improve the robustness of representation engineering-based defensive methods by reducing the attack success rate (ASR) of the Llama 3 8B model from 29% to 5% against embedding attacks across diverse configurations, and from 14% to 0% against the REINFORCE-GCG (Geisler et al., 2025a) input space attack.

Our main contributions are as follows:

- 1. We propose a novel method for LLM safety based on contrastive representation learning, generalizing and improving upon existing methods such as circuit breakers (Zou et al., 2024) and RepBend (Yousefpour et al., 2025).
- 2. We introduce a new training strategy based on representation adversarial training (Casper

et al., 2024) that improves the sampling procedure of harmful representations, further increasing the robustness of our proposed method against embedding attacks.

3. We provide experimental validation of our method, demonstrating gains in robustness against both input space and embedding space attacks without sacrificing the model's general capabilities.

2 Problem formulation

Large Language Models (LLMs) are deep neural networks, typically based on the Transformer architecture (Vaswani et al., 2017), that have been trained at a massive scale on human textual data. The computation of a token y_t given previous tokens $y_{< t}$ can be described as follows:

$$\mathbf{h}_0 = \operatorname{Emb}(y_{< t}) \tag{1}$$

$$\mathbf{h}_l = T_l(\mathbf{h}_{l-1}) + \mathbf{h}_{l-1} \tag{2}$$

$$\mathbf{o} = \mathrm{FFN}(\mathbf{h}_L) \tag{3}$$

where $\mathbf{h}_0 \in \mathbb{R}^{t-1 \times d}$ is the sequence of input embeddings for the previous tokens, T_l is the transformer block at layer l, $\mathbf{h}_l \in \mathbb{R}^{t-1 \times d}$ is the output of the l-th layer, FFN is a final feed-forward layer, and $\mathbf{o} \in \mathbb{R}^{t-1 \times |\Sigma|}$ represent the output log probabilities over the vocabulary Σ for each position. Finally, y_t is sampled from the output distribution \mathbf{o} .

For a prompt x and a response y, we write $\mathbf{h}_l(x,y)$ as the inner representation of the model at layer l. We consider a prompt and response pair $(x,y)_h$ to be harmful if its content violates the developer rules of the model. Conversely, a pair $(x,y)_b$ benign does not violate developer rules.

Objective Our study specifically targets prompts that encourage illegal, immoral, unethical, or dangerous conduct and actions. The objective is to minimize the probability that a model outputs a harmful reply y under any token or embedding input x.

We use the same proxy objective as the circuitbreaking method of working in the representation level, as while input and output formats can change between different generations of the harmful behavior, the representation space symbolizes the same concept under different formats. We aim to create a new safer model that "breaks" when given a harmful prompt, and works as expected when given a benign prompt.

3 Related work

3.1 AI Safety

Input-Space Jailbreaking Attacks A jailbreak is a prompt specifically designed to bypass the model's safety mechanisms and elicit a harmful response, and can be broadly categorized into token-level jailbreaks which optimize a harmful adversarial sequence of tokens appended to the prompt (e.g. Greedy Coordinate Gradient (Zou et al., 2023)), or prompt-level jailbreak which optimize the entire prompt into human-readable jailbreak prompts (e.g Prompt Automatic Iterative Refinement (PAIR) (Chao et al., 2024)).

Embedding-Space Jailbreaking Attacks Embedding attacks directly manipulate the output of the model's embedding layer to produce a harmful response (Schwinn et al., 2025). For a prompt x, instead of optimizing an adversarial suffix $y \in \Sigma^t$, the attacker can optimize an embedding $e \in \mathbb{R}^{t \times d}$ that is appended to the prompt embedding $\mathrm{Emb}(x)$, to force the model to start with a positive reply. Mathematically, we can write $\mathbf{h}_0 = \mathrm{Emb}(x) \oplus e$ where \oplus is the concatenation operator. Embedding attacks are more powerful than input space attacks as they operate on a lower level and have access to the entire embedding space of the model.

Jailbreak Defenses Jailbreak defenses aim to prevent or mitigate the effects of jailbreak attacks on LLMs. Current defenses fall into two main categories: Pipeline defenses external to the model, which do not modify the model weights but rather add components such as filters (Jain et al., 2023), textual transformations (Robey et al., 2024; Yuan et al., 2024b) or guardrail models (Zeng et al., 2024b), and internal defenses which alter the model itself, by fine-tuning on preference data (Kaufmann et al., 2024) or editing problematic layers (Zhao et al., 2024).

More details on jailbreaking attacks and defenses are found in Appendix A.1. In this work, we focus on a subset of internal defenses that manipulate the model's internal representations to increase robustness against jailbreak attacks.

Internal Defenses based on Representation Engineering Representation Engineering (Zou et al., 2025) focuses on internal representational spaces to understand and improve the behavior of LLM, and is an alternative to mechanistic interpretability. Circuit breakers (Zou et al., 2024) and RepBend

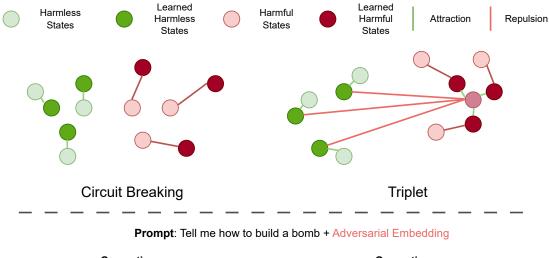


Figure 1: Comparison of the Triplet defense with the Circuit Breaking defense. Contrary to other adversarial defense methods, circuit breaking aims to break generation at harmful content instead of refusing to answer harmful tasks. It fine-tunes models to keep learned harmless states (or representations) close together while separating newly learned harmful states from their original counterparts, without additional constraints. In contrast, the Triplet defense additionally pulls learned harmful states together and pushes them away from learned harmless states, which increases contrast and robustness to embedding-space attacks.

(Yousefpour et al., 2025) are the two main representation engineering-based defenses. These methods share common principles:

- 1. They manipulate the inner representations \mathbf{h}_l of the model, rather than only the inputs and outputs.
- 2. They define training loss functions over these inner representations to guide and optimize model behavior.
- 3. They use datasets of prompts and responses labeled as "benign" or "harmful" to establish "good" and "bad" inner representations.

The training loss of circuit breaking focuses on preserving benign representations and pushing the new harmful representations far from previous harmful representations using a cosine similarity loss, and is defined as

$$\mathcal{L}_{CB} = \alpha \cdot \left\| \mathbf{h}_b - \mathbf{h}_b' \right\|_2^2 + \beta \cdot \text{ReLU}(\cos_{-}\text{sim}(\mathbf{h}_h, \mathbf{h}_h'))$$
 (4)

where \mathbf{h} is the inner representation of the initial fixed model at layer l, \mathbf{h}' is the new inner representation of the model with circuit breaking. α

Circuit breaking is an effective defense against a wide range of input-space attacks, and decreases the quality of successful attacks. For instance, while the REINFORCE adversarial attack (Geisler et al., 2025a) achieves a high attack success rate (ASR) with the HarmBench judge (Mazeika et al., 2024), the resulting responses often display stutter-like behavior in later stages of generation, rendering them mostly unusable. However, circuit breaking is less effective at defending against embedding-level jailbreaking attacks (Schwinn and Geisler, 2024).

Representation Bending (RepBend) (Yousefpour et al., 2025) extends circuit breaking by replacing the cosine similarity-based distance with an L2 distance, and introduces an additional regularization term to enforce similarity among harmful representations. However, training this method is challenging, as the distance terms can diverge to negative infinity, requiring careful stopping criteria and making it difficult to maintain a balance of hyperparameters.

3.2 Contrastive Representation Learning

The objective of contrastive representation learning (CRL) is to train models to produce a representational space in which similar (positive) inputs are mapped close to each other, while dissimilar (negative) inputs are mapped far apart. Rather than solely relying on labeled data, contrastive representation can learn meaningful representations by leveraging the inherent structure of the data itself. It has achieved notable success in a variety of fields, such as computer vision (Le-Khac et al., 2020; Schroff et al., 2015), natural language processing (Mikolov et al., 2013; Rim et al., 2021), and multi-modal learning (Radford et al., 2021).

The triplet loss (Schroff et al., 2015) is a popular loss function used in contrastive learning, originally developed for image embeddings and face recognition and later adopted for text embeddings (Reimers and Gurevych, 2019).

It encourages an anchor data point a to be closer to a positive sample p (similar to the anchor) than to a negative sample n (dissimilar to the anchor) by at least a margin m:

$$\mathcal{L}_T = \text{ReLU}(d(a, p) - d(a, n) + m)$$
 (5)

where d(.,.) is a distance metric, typically the Euclidean distance or cosine distance, and $\operatorname{ReLU}(x) = \max(0,x)$ is the rectified linear unit function. A visualization of the triplet loss objective is shown in Appendix A.1.

4 Method

We first describe the desirable properties of our new, more robust representation space. Our proxy objective is to construct a new representation space \mathbf{h}' that induces the following properties, for an index i and benign and harmful representations $\mathbf{h}'_{b,i}$ and $\mathbf{h}'_{b,i}$:

- 1. $\mathbf{h}'_{b,i}$ and $\mathbf{h}_{b,i}$ should be similar to each other, as the new model is expected to maintain similar behavior to the original model for benign use cases. Exact matching is not required for all representations, but top logits of benign behaviors should closely match.
- 2. $\mathbf{h}'_{h,i}$ and $\mathbf{h}_{h,i}$ should be dissimilar to each other, because if the general structure of the representation space is preserved, the old harmful representations maintain their harmful nature in the new representation space.

- 3. $\mathbf{h}'_{b,i}$ and $\mathbf{h}'_{h,i}$ should be dissimilar, allowing the model to distinguish between benign and harmful representations.
- 4. $\mathbf{h}'_{h,i}$ and $\mathbf{h}'_{h,j}$ should be similar to each other, preventing the model from generating fine-grained responses to harmful queries and promoting the generation of uniform replies such as refusals or warnings.

Interpreting the circuit breaking loss The circuit breaking loss described in Equation 4 can be interpreted as a contrastive loss, similar in spirit to the DrLIM loss (Hadsell et al., 2006), which is one of the first contrastive losses.

For input vectors \mathbf{X}_1 and \mathbf{X}_2 belonging to a class $Y \in \{0, 1\}$ the DrLIM loss is defined as

$$\mathcal{L}_{DrLIM} = (1 - Y)\frac{1}{2} \|\mathbf{X}_1 - \mathbf{X}_2\|_2^2 + (Y)\frac{1}{2} \max(0, m - d(\mathbf{X}_1, \mathbf{X}_2))$$
(6)

where $d(\mathbf{X}_1, \mathbf{X}_2)$ is a distance and m is a margin hyperparameter.

This objective reduces to the circuit breaking objective when using the distance $d(\mathbf{X}_1, \mathbf{X}_2) = 1 - \cos_{\sin}(\mathbf{X}_1, \mathbf{X}_2)$, a hard margin m = 1, harmfulness labels Y, and sampling $\mathbf{X}_1, \mathbf{X}_2$ from the original and fine-tuned models.

In CRL tasks, the DrLIM loss has been largely supplanted by more effective objectives, notably the Triplet loss and the InfoNCE loss (van den Oord et al., 2019), as they are more flexible and induce greater contrasts between the representations. Motivated by these advances, we use a triplet loss to learn a robust representation space for LLM defense.

4.1 Our Triplet-Based Loss

Taking inspiration from the circuit breaking loss function, we propose a general alternative loss function that fits all wanted properties.

Let $d_{h,p}$, $d_{h,n}$, $d_{b,p}$ and $d_{b,n}$ be distance functions on representations, and i an index. We define a harmful triplet loss as:

$$\mathcal{L}_{triplet}(h_i) = \text{ReLU}(d_{hp}(\mathbf{h}'_{h,i}, \mathbf{p}_{h,i}) - d_{hn}(\mathbf{h}'_{h,i}, \mathbf{h}_{h,i}) + m_h) \quad (7)$$

This loss encourages new harmful representations $\mathbf{h}'_{h,i}$ to be distant from the old harmful representations $\mathbf{h}_{h,i}$, and close to some positive $\mathbf{p}_{h,i}$. In contrast to circuit breaking and RepBend, our

approach focuses on relative rather than absolute distances between representations, as relative distances are more meaningful in embedding spaces. Conversely, we define a benign triplet loss as:

$$\mathcal{L}_{triplet}(b_i) = \text{ReLU}(d_{bp}(\mathbf{h}_{b,i}, \mathbf{h}'_{b,i}) - d_{bn}(\mathbf{h}'_{b,i}, \mathbf{n}_{b,i}) + m_b)$$
(8)

This loss encourages new benign representations $\mathbf{h}'_{b,i}$ to be close to the old benign representations $\mathbf{h}_{b,i}$, and far from some negative $\mathbf{n}_{b,i}$. We write the final, unified triplet loss as a weighted sum of the two triplet losses:

$$\mathcal{L}_{triplet} := \alpha \mathcal{L}_{triplet}(b_i) + \beta \mathcal{L}_{triplet}(h_i)$$
 (9)

with hyperparameters α and β controlling the importance of the losses.

Property	RepBend	CB	Triplet
$\mathbf{h}_{b,i}^{\prime}pprox\mathbf{h}_{b,i}$	✓	√	✓
$\mathbf{h}_{h,i}^{'} otseph\mathbf{h}_{h,i}$	\checkmark	\checkmark	\checkmark
$\mathbf{h}_{h,i}^{\prime} otsete{\mathbf{h}_{h,i}} \mathbf{h}_{b,i}^{\prime} otsete{\mathbf{h}_{h,i}} $	X	X	\checkmark
$\mathbf{h}_{h,i}^{\prime\prime}pprox\mathbf{h}_{h,j}^{\prime\prime}$	\checkmark	X	\checkmark

Table 1: Safety representation engineering methods and their properties. Properties are expressed in terms of similarity (\approx) or dissimilarity $(\not\approx)$

Relation to Circuit Breakers and RepBend We demonstrate that both the circuit breaking loss and the RepBend loss are simplified special cases of our triplet loss. The full derivation and proofs are in Appendix A.10. Table 1 outlines the key differences between the three losses. Specifically, the circuit breaking loss lacks mechanisms for separating benign representations from harmful ones, and for clustering the harmful representations. The RepBend loss focuses on clustering harmful representations, but does not explicitly separate benign representations from harmful ones. Our triplet loss formulation generalizes both methods by incorporating these properties and allowing any valid

Choice of the Positive and Negative Samples

pseudodistances.

The choice of $\mathbf{p}_{h,i}$ and $\mathbf{n}_{b,i}$ is important, as these samples will guide the new representations to new, better directions. In this work, we use the mean of the new harmful representations as a positive sample for the harmful triplet loss, and as a negative sample for the benign triplet loss, thereby strengthening the separation between the two classes of

representations. A visualization of our method compared to circuit breakers is shown in Figure 1. Future work should explore the use of other choices of $\mathbf{p}_{h,i}$ and $\mathbf{n}_{b,i}$.

Final Triplet Loss Let α , β and γ be hyperparameters controlling the importance of the losses. Let \mathbf{d}_{bp} , \mathbf{d}_{hn} , \mathbf{d}_{hp} and \mathbf{d}_{hn} be distances chosen by the user. Let $\mathbf{h}_{b,i}$ and $\mathbf{h}_{h,i}$ be the benign and harmful representations for a batch with N different benign and harmful prompts. Let $\hat{\mathbf{h}}' = \frac{1}{N} \sum_{h=1}^{N} \mathbf{h}'_h$ be the mean of the harmful representations for a batch. Let \mathbf{D}_{KL} the Kullback-Leibler divergence on benign model logits between the new and the original model. Our final triplet loss is defined as:

$$\mathcal{L}_{Triplet} = \alpha \cdot \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}_{triplet}(b_i)$$
 (10)

$$+ \beta \cdot \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}_{triplet}(h_i)$$
 (11)

$$+ \gamma \cdot D_{KL}(\mathcal{M}(x_b) \parallel \mathcal{M}'(x_b))$$
 (12)

Algorithm 1 describes the training procedure for the triplet model defense. The model weights are optimized until convergence on batches of benign and harmful prompt-response pairs.

4.2 Combining Representation Engineering with Adversarial Training

Most LLMs are shipped with built-in safety features that prevent them from outputting harmful responses to plain harmful queries. As such, gathering harmful representations $\mathbf{h}'_{h,i}$ from plain queries can lead to representations that are not truly informative of dangerous model behavior. Inspired by work in hard negative mining, which focuses on learning on challenging negative samples (Robinson et al., 2021), we propose to address this issue by integrating adversarial training in the representation space, by explicitly generating "hard" harmful representations via attacks.

Adversarial Hard Negative Mining In contrastive learning, "hard negatives" are challenging negative examples that are easily confused with positive examples. For safety representation engineering, we define hard negatives as harmful representations that closely resemble benign ones. Rather than relying on plain harmful representations, we make use of adversarial hard negative mining (Hughes et al., 2018).

Algorithm 1 Triplet Model Defense

Require: Frozen original model \mathcal{M} ; Trainable defense model \mathcal{M}' , Benign dataset \mathcal{D}_b , harmful dataset \mathcal{D}_h ; Number of steps T; batch size N; Hyperparameters $\alpha, \beta, \gamma, m_b, m_h$

- 1: **for** t = 1, ..., T **do**
- 2: Sample a batch $x_b \sim \mathcal{D}_b$, $x_h \sim \mathcal{D}_h$
- Compute original representations $\mathbf{h}_{b,i}, \mathbf{h}_{h,i}$ using \mathcal{M} 3:
- Compute new representations $\mathbf{h}'_{b,i}$, $\mathbf{h}'_{h,i}$ using \mathcal{M}'
- Compute $\hat{\mathbf{h}}' = \frac{1}{N} \sum_{i=1}^{N} \mathbf{h}'_{h,i}$ 5:
- $\mathcal{L}_{\text{benign}} = \frac{1}{N} \sum_{i=1}^{N} \max \left(0, d_{bp}(\mathbf{h}_{b,i}, \mathbf{h}'_{b,i}) d_{bn}(\mathbf{h}'_{b,i}, \hat{\mathbf{h}'}) + m_b \right)$ 6:
- $\mathcal{L}_{\text{harmful}} = \frac{1}{N} \sum_{i=1}^{N} \max \left(0, d_{hp}(\mathbf{h}'_{h,i}, \hat{\mathbf{h}'}) d_{hn}(\mathbf{h}'_{h,i}, \mathbf{h}_{h,i}) + m_h \right)$ 7:
- $\mathcal{L}_{KL} = D_{KL}(\mathcal{M}(x_b) \parallel \mathcal{M}'(x_b))$ 8:
- $\begin{aligned} \mathcal{L}_{\text{Triplet}} &= \alpha \cdot \mathcal{L}_{\text{benign}} + \beta \cdot \mathcal{L}_{\text{harmful}} + \gamma \cdot \mathcal{L}_{\text{KL}} \\ \text{Update parameters of } \mathcal{M}' \text{ using } \mathcal{L}_{\text{Triplet}} \end{aligned}$ 9:
- 10:
- 11: end for

An adversarial attack neural network module Attack_l is introduced at a randomly selected layer l. The module is inserted between two transformer blocks in the residual stream, and is trained using a Negative Log Likelihood (NNL) loss on harmful responses. The model is active when sampling new harmful representations, and is periodically retained as model parameters are updated. This module finds adversarial hard negatives, allowing the defense to counteract a wider diversity of harmful representations. Appendix A.2 contains more details on the training of our adversarial attack modules.

In summary, our complete method views safety representation engineering as a contrastive learning problem, optimizes a triplet-based loss function, and combines it with adversarial training to defend against harmful prompts more robustly.

Experimental setup

Models We evaluate our method on two widely used open-source models: Llama 3 8B Instruct (Grattafiori et al., 2024) and Mistral 7B Instruct v0.2 (Jiang et al., 2023). These models are standard baselines for adversarial defense in the literature, which enables direct comparison with prior work. Additional results on two extra models are provided in Appendix A.6.

Datasets We adapt the training pipeline of Yousefpour et al. (2025) for our defense method. For benign data, we use UltraChat (Ding et al., 2023), a large-scale dataset of over 1.5 million multi-turn dialogues that cover a wide range of topics such as art, history, literature, politics and technology. For harmful data, we use WildGuardMix (Han et al., 2024) which contains a broad spectrum of jailbroken prompts and responses, and Wild-Jailbreaks (Jiang et al., 2024), a synthetic dataset of harmful prompt-response pairs, featuring both straightforward and complex jailbreak prompts. From these sources, we randomly select 10'000 benign and 10'000 harmful samples to construct a balanced training set.

Details on hyperparameter choices can be found in Appendix A.3.

Experiments

In this section, we demonstrate the increased robustness of our methods against various attacks compared to existing methods and explore three research questions (RQ) related to robustness and general performance.

Overall defensive performance

In this section, we compare the performance of our approach against different models and defenses.

Method We use the publicly available defensive models created by the authors of circuit breakers (Zou et al., 2024) and RepBend (Yousefpour et al., 2025). We evaluate our defense using the Harm-Bench safety benchmark (Mazeika et al., 2024). For embedding attacks, we adopt a variant of the attack described by Zou et al. (2024). To ensure robustness across hyperparameter choices, we use 6 different hyperparameter configurations. For each behavior, we select the result from the configuration that produces the most harmful response out of the six different runs. For input-space attacks,

	REINFORCE-GCG				GCG			Embedding		
Defense	HB	SR	Score	HB	SR	Score	НВ	SR	Score	
Original model	52.50	40.00	42.87	31.25	18.75	23.66	100.00	90.24	81.89	
Circuit breakers	13.75	3.75	9.50	2.86	1.43	4.25	90.24	29.27	30.61	
RepBend	11.25	6.25	11.27	2.86	0.00	1.65	73.17	39.02	39.00	
Triplet	0.00	0.00	0.48	0.00	0.00	0.43	65.85	12.20	14.57	
Triplet + Adv	3.75	2.50	6.99	0.00	0.00	1.36	75.61	4.88	8.70	

Table 2: Attack success rates (ASR) using HarmBench (HB) and StrongREJECT (SR) across attack types, for various defenses (Llama 3 8B Instruct). For GCG and REINFORCE-GCG attacks, each behavior was tested on a single attempt evaluated over 80 Behaviors. For embedding attacks, results were computed over 41 behaviors, with six attempts per behavior using different hyperparameter configurations (246 runs per model). The best result for each behavior was used. StrongREJECT scores are reported on a 0–100 scale.

we use GCG (Zou et al., 2023) and REINFORCE-GCG (Geisler et al., 2025a) with base configurations. To assess harmfulness, we use the binary HarmBench judge to get adversarial success rates (HB ASR) and the fine-grained StrongREJECT (Souly et al., 2024) fine-tuned classifier to get harmfulness scores (SR Score) and adversarial success rates (SR ASR) for scores above 0.5. Full details of attack configurations and evaluation settings are provided in Appendix A.3

Results Tables 2 report the ASRs for embedding and GCG attacks on the Llama 3 8B model. All defense methods achieve substantial improvements over the base model, which shows ASRs above 90% for embedding attacks. Among the defenses evaluated, our triplet defenses outperform both circuit breakers and RepBend. In particular, the triplet defense achieves ASRs of 0% for both RE-INFORCE and GCG, while the triplet defense with adversarial hard negative mining achieves the lowest embedding SR ASR of 4.88% and the lowest harmfulness score of 8.70.

Figure 2 shows embedding attack success rates for Llama 3 8B with two additional adversarial defenses: Refusal Feature Ablation Training (ReFAT) (Yu et al., 2025) and Latent Adversarial Training (LAT) (Sheshadri et al., 2024). The Triplet model substantially improves embedding ASRs compared to these baselines. Results for the Mistral 7B model are shown in Appendix A.6. Although the reported ASRs are higher than for the Llama model, the triplet defense outperforms the evaluated baselines.

Throughout our experiments, we find that the HarmBench classifier consistently produces higher ASRs than the StrongREJECT classifier. This is likely due to the tendency of HarmBench to classify responses as harmful responses based on the

initial response tokens, even if the rest of the response is nonsensical or incoherent. Appendix Table 21 shows an example of a generation classified as harmful by the HarmBench ASR, despite being practically harmless. These findings highlight how differences in harmfulness evaluation criteria can substantially influence measured outcomes.

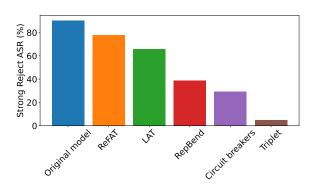


Figure 2: Embedding Attack success rate (ASR) using StrongREJECT for various defenses (Llama 3 8B Instruct.

6.2 RQ1: How robust is our defense to different attack configuration choices?

The goal of this research question is to evaluate the robustness of our defense to different embedding configurations.

Method We compare the ASRs of embedding attacks across the six different hyperparameter configurations, with full details provided in Appendix A.3. For each defense, we report the best, worst, and mean ASR obtained over all configurations.

Results Table 3 presents our results. Both Rep-Bend and circuit breakers exhibit significant variability across attack hyperparameter configurations.

Defense	HarmBench ASR			Stroi	StrongREJECT ASR			StrongREJECT Score		
	mean	min	max	mean	min	max	mean	min	max	
Original model	77.33	54.00 ¹	98.00^{3}	48.16	26.10^4	63.34 ³	53.25	24.39^4	73.17 ⁵	
RepBend	24.50	2.00^{5}	37.00^2	10.36	4.00^{5}	22.06^2	8.54	2.44^{0}	19.51^2	
Circuit breakers	38.67	27.00^2	54.00^{1}	9.32	3.41^{5}	14.53 ⁰	6.91	0.00^{5}	12.20^{0}	
Triplet	23.83	17.00^{3}	32.00^{1}	3.55	1.16^{4}	9.46^{2}	2.44	0.00^{0}	9.76^{2}	
Triplet + Adv	24.40	10.00 ¹	41.00^2	2.23	1.10^{3}	4.28°	0.49	0.00^{1}	2.44^{0}	

Table 3: Mean, maximum, and minimum embedding attack ASR across six different hyperparameter configurations (Llama 3 8B). Colored superscript indicates the configuration index for which the ASRs were achieved.

In particular, using configuration 2, we get a StrongREJECT ASR of 20% on RepBend and 2% on circuit breakers, while configuration 0 results in an ASR of 2% and 12% respectively. In comparison, our triplet defense consistently demonstrates low ASRs, with a worst-case StrongREJECT ASR of 2%. The complete results are provided in Appendix Figure 19. These findings highlight the necessity of evaluating defenses across diverse attack configurations to accurately assess their robustness, as relying on a single configuration could bias results in favor of a particular defense.

6.3 RQ2: Does applying our defense affect the general performance of the model?

The objective of this research question is to determine whether the application of our defensive mechanism affects the general language capabilities of the model.

Method We assess the general performance of our trained models on a suite of benchmarks, including MMLU, HellaSwag, TruthfulQA, and GSM8K. See Appendix A.3 for more details.

Results As shown in Table 4, our triplet method achieves a performance comparable to the base model, which indicates that our approach preserved the model's general capabilities. Notably, the defenses showcase an improvement on the TruthfulQA (MC) benchmark, likely because their increased tendency to reject harmful responses also leads them to reject untruthful content, which is often harmful.

In contrast, the performance of the RepBend model decreases significantly on the generation-based benchmarks GSM8K and TruthfulQA (Gen), with accuracy dropping from 75% to 49% on GMS8k, and from 46% and 2% on TruthfulQA (Gen). This suggests overfitting to the defense objective at the expense of general language performance. Examples of generations of GSM8K

for RepBend can be seen in Appendix Table 21, in which the model fails to answer the questions. In contrast, our trained triplet models do not suffer from this issue, maintaining both GSM8K and TruthfulQA (Gen) performance close to the base model. These findings highlight the strengths of our approach and illustrate that defenses can have unintended side effects, which may themselves have harmful consequences in sensitive applications.

6.4 RQ3: How does our defense generalize to out-of-distribution input/output formats?

While adversarial attack success rates are valuable for measuring robustness in plain text, they do not fully capture a defense's capability to generalize to out-of-distribution response formats. To address this, we introduce a new evaluation metric for safety representation engineering defenses, based on relative distances.

Method We apply random capitalization (following the Best-of-N jailbreak attack protocol (Hughes et al., 2024)) and translations to five languages to 159 behaviors of the HarmBench benchmark, generating a set of augmented prompts and responses $\mathcal{A}(b)$ for each behavior b. For each defense, we compute the Mean Minimum Distance Ratio (MMDR) as the average over behaviors of the smallest ratios between the distance of augmented and original harmful representations:

$$MMDR_d = \frac{1}{|N|} \sum_{i \in N} \min_{a \in \mathcal{A}(h_i)} \frac{d(\mathbf{h}'_a, \mathbf{h}_a)}{d(\mathbf{h}'_{h,i}, \mathbf{h}_{h,i})} \quad (13)$$

Averaged over all model layers, MMDR quantifies the model's worst-case generalization to out-of-distribution augmentations.

Results As shown in Table 5, the full triplet defenses achieve higher MMDR values, up to 0.8 for both distance metrics, compared to the circuit

-	ARC (Easy)	GSM8K	HellaSwag	MMLU	TruthfulQA		
					Gen	MC1	MC2
Original model	81.61	75.36	57.75	63.72	46.39	36.23	51.67
Circuit breakers	81.44 (-0.17)	75.44 (+0.08)	57.46 (-0.29)	63.57 (-0.15)	48.23 (+1.84)	36.96 (+0.73)	51.61 (-0.05)
RepBend	80.98 (-0.63)	49.05 (-26.31)	60.58 (+2.83)	60.26 (-3.46)	2.08 (-44.31)	41.00 (+4.77)	60.05 (+8.38)
Triplet	81.27 (-0.34)	74.30 (-1.06)	59.62 (+1.87)	63.85 (+0.13)	45.65 (-0.73)	40.76 (+4.53)	55.37 (+3.70)
Triplet + Adv	81.99 (+0.38)	74.91 (-0.45)	60.70 (+2.95)	63.38 (-0.34)	44.55 (-1.84)	42.96 (+6.73)	57.29 (+5.63)

Table 4: Performance comparison of models on general capability benchmarks (Llama 3 8B Instruct).

Distance	$\overline{ ext{MMDR}_{d_2}}$	$\mathbf{MMDR}_{d_{\mathrm{cos}}}$
Circuit breakers	0.63	0.49
Triplet A1: CB	0.70	0.54
RepBend	0.70	0.64
Triplet A2: RepBend	0.78	0.64
Triplet A3: Full	0.80	0.66
Triplet A4: Full + Adv	0.80	0.70

Table 5: Generalization of the defenses to different data augmentations (Llama 3 8B Instruct) with the L2 norm d_2 and the cosine distance d_{\cos}

breaking (0.6) and RepBend (0.7) defenses. A value close to 1 indicates strong generalization, meaning the defense modifies augmented harmful representations similarly to the unaugmented ones, while a value close to 0 would imply that some augmentations bypass the defense. Notably, the MMDR increases as more loss terms are incorporated, highlighting the effectiveness of our approach in generalizing to different input and output formats.

6.5 Ablation study

To study the impact of our loss functions, we conduct an ablation study in which our defense methods are trained with different loss configurations.

Specifically, we train a model A1 by removing the $d_{b,n}$ and $d_{h,p}$ components, making it closely related to circuit breaking. For model A2, we ablate only $d_{b,n}$, resulting in a formulation similar to RepBend but with a margin-based objective. Model A3 retains all loss terms, while model A4 incorporates adversarial hard negative mining and all loss terms. Details and full results are provided in Appendix A.4

Our results show that A1 performs similarly to circuit breaking, validating the correspondence between the loss formulations. A2 consistently outperforms both A1 and RepBend, highlighting the advantage of our margin-based triplet learning objective. A3 and A2 perform comparably overall, with A3 achieving better results in input-space at-

tacks and A2 on embedding space attacks. A4 achieves the best performance. Table 5 further shows that for our studied models, removing loss terms leads to a decrease in MMDR. These findings demonstrate the importance of the $d_{b,n}$ term in our triplet loss, as well as the additional benefit of adversarial hard negative mining.

7 Conclusion

This work presents a novel method for improving the robustness of LLMs against adversarial attacks based on contrastive representation learning and adversarial hard negative mining. Our findings demonstrate notable robustness improvement while maintaining the model's general capabilities.

Ethics Statement

This section discusses the ethical considerations that arise from the development and deployment of defensive methods for AI models. First, the development of defenses may lead to overconfidence in the safety of AI models, which in turn could encourage the deployment of less safe models. In addition, better defenses may lead to the development of stronger, more sophisticated attacks, thereby increasing the risk of misuse for widespread AI models with fewer security measures. Lastly, our representation space attack module used for adversarial hard negative mining could be misused by malicious actors to circumvent the defenses of AI models in a white-box setting. Nevertheless, these potential downsides are outweighed by the benefits of developing better defenses. These concerns underscore the need for responsible use and deployment of research findings in the field of AI safety.

Limitations

Despite the promising results of our method, several limitations should be acknowledged. First, while robustness improves significantly on the

Llama model, the Mistral model remains more vulnerable to attacks, even though our method still outperforms circuit breakers. This suggests that further tuning or architecture-specific adaptations may be required to achieve a strong robustness across models.

Secondly, given the computationally intensive nature of adversarial training and jailbreak generation, the choice of hyperparameters and training strategy may not be optimal. Furthermore, due to the 2-GPU-hour cost per REINFORCE attack, our method was evaluated on 80 HarmBench behaviors in the validation set compared to the 300 behaviors in the training set. Although we expect similar trends to hold, this assumption has not been empirically verified.

Third, while our method is robust to a variety of attacks, it is not guaranteed to be robust to all, especially to attacks in the representation space. Using more attack configurations and attempts per behavior may also lead to high attack success rates, at the cost of increased computation time.

Finally, like other representation engineering-based methods, our trained models may result in incoherent and ineligible behavior if the model misinterprets benign inputs as harmful, which in turn could lead to harmful consequences in some critical settings. Therefore, careful considerations and additional safeguards may be necessary before deploying these methods in real-world applications. Finally, our method with adversarial hard negative mining requires a moderate training time of up to 12 hours on a single H100 GPU for the Llama 3 8B model. This requirement may limit the scalability of our approach to much larger models with hundreds of billions of parameters.

Acknowledgements

We thank Kellin Pelrine, Roger Grosse, and Stephen Casper for their feedback on our work. This material is based in part upon work supported by the German Federal Ministry of Education and Research (BMBF): Tübingen AI Center, FKZ: 01IS18039B; by the Machine Learning Cluster of Excellence, EXC number 2064/1 – Project number 390727645; by Schmidt Sciences SAFE-AI Grant; by NSERC Discovery Grant RGPIN-2025-06491; by a National Science Foundation award (#2306372); by a Swiss National Science Foundation award (#201009) and a Responsible AI grant by the Haslerstiftung; as part of the "Swiss

AI initiative" by a grant from the Swiss National Supercomputing Centre (CSCS) under project IDs a07 and a130 on Alps. Additional calculations were carried out on the ETH Euler cluster and the Center for AI Safety cluster.

References

Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, and 110 others. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *Preprint*, arXiv:2404.14219.

Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. 2025. Jailbreaking leading safety-aligned llms with simple adaptive attacks. *Preprint*, arXiv:2404.02151.

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

Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal Valko, and Rémi Munos. 2023. A general theoretical paradigm to understand learning from human preferences. *Preprint*, arXiv:2310.12036.

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, and 12 others. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *Preprint*, arXiv:2204.05862.

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

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

Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric Wong. 2024. Jailbreaking black box large language models in twenty queries. *Preprint*, arXiv:2310.08419.

Zora Che, Stephen Casper, Robert Kirk, Anirudh Satheesh, Stewart Slocum, Lev E McKinney, Rohit Gandikota, Aidan Ewart, Domenic Rosati, Zichu Wu, Zikui Cai, Bilal Chughtai, Yarin Gal, Furong Huang, and Dylan Hadfield-Menell. 2025. Model

- tampering attacks enable more rigorous evaluations of llm capabilities. *Preprint*, arXiv:2502.05209.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. *Preprint*, arXiv:2002.05709.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *Preprint*, arXiv:1803.05457.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *Preprint*, arXiv:2110.14168.
- Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. 2023. Toxicity in chatgpt: Analyzing persona-assigned language models. *Preprint*, arXiv:2304.05335.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. *Preprint*, arXiv:2305.14233.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, and 5 others. 2024. The language model evaluation harness.
- Jonas Geiping, Alex Stein, Manli Shu, Khalid Saifullah, Yuxin Wen, and Tom Goldstein. 2024. Coercing llms to do and reveal (almost) anything. *Preprint*, arXiv:2402.14020.
- Simon Geisler, Tom Wollschläger, M. H. I. Abdalla, Vincent Cohen-Addad, Johannes Gasteiger, and Stephan Günnemann. 2025a. Reinforce adversarial attacks on large language models: An adaptive, distributional, and semantic objective. *Preprint*, arXiv:2502.17254.
- Simon Geisler, Tom Wollschläger, M. H. I. Abdalla, Johannes Gasteiger, and Stephan Günnemann. 2025b. Attacking large language models with projected gradient descent. *Preprint*, arXiv:2402.09154.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

- Raia Hadsell, Sumit Chopra, and Yann Lecun. 2006. Dimensionality reduction by learning an invariant mapping. pages 1735 1742.
- Seungju Han, Kavel Rao, Allyson Ettinger, Liwei Jiang, Bill Yuchen Lin, Nathan Lambert, Yejin Choi, and Nouha Dziri. 2024. Wildguard: Open one-stop moderation tools for safety risks, jailbreaks, and refusals of llms. *Preprint*, arXiv:2406.18495.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Preprint*, arXiv:2009.03300.
- Xinyi Hou, Yanjie Zhao, Yue Liu, Zhou Yang, Kailong Wang, Li Li, Xiapu Luo, David Lo, John Grundy, and Haoyu Wang. 2024. Large language models for software engineering: A systematic literature review. *Preprint*, arXiv:2308.10620.
- Brian R. Y. Huang, Maximilian Li, and Leonard Tang. 2024. Endless jailbreaks with bijection learning. *Preprint*, arXiv:2410.01294.
- John Hughes, Sara Price, Aengus Lynch, Rylan Schaeffer, Fazl Barez, Sanmi Koyejo, Henry Sleight, Erik Jones, Ethan Perez, and Mrinank Sharma. 2024. Bestof-n jailbreaking. *Preprint*, arXiv:2412.03556.
- Lloyd H. Hughes, Michael Schmitt, and Xiao Xiang Zhu. 2018. Generative adversarial networks for hard negative mining in cnn-based sar-optical image matching. In *IGARSS 2018 2018 IEEE International Geoscience and Remote Sensing Symposium*, pages 4391–4394.
- Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping yeh Chiang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. 2023. Baseline defenses for adversarial attacks against aligned language models. *Preprint*, arXiv:2309.00614.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.
- Liwei Jiang, Kavel Rao, Seungju Han, Allyson Ettinger, Faeze Brahman, Sachin Kumar, Niloofar Mireshghallah, Ximing Lu, Maarten Sap, Yejin Choi, and Nouha Dziri. 2024. Wildteaming at scale: From in-the-wild jailbreaks to (adversarially) safer language models. *Preprint*, arXiv:2406.18510.
- Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. 2024. A survey of reinforcement learning from human feedback. *Preprint*, arXiv:2312.14925.

- Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2021. Supervised contrastive learning. *Preprint*, arXiv:2004.11362.
- Diederik P. Kingma and Jimmy Ba. 2017. Adam: A method for stochastic optimization. *Preprint*, arXiv:1412.6980.
- Phuc H. Le-Khac, Graham Healy, and Alan F. Smeaton. 2020. Contrastive representation learning: A framework and review. *IEEE Access*, 8:193907–193934.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. *Preprint*, arXiv:2109.07958.
- 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. 2024. Harmbench: A standardized evaluation framework for automated red teaming and robust refusal. *Preprint*, arXiv:2402.04249.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *Preprint*, arXiv:1301.3781.
- Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: A survey. *Preprint*, arXiv:2402.06196.
- Mohamed Nejjar, Luca Zacharias, Fabian Stiehle, and Ingo Weber. 2024. Llms for science: Usage for code generation and data analysis. *Preprint*, arXiv:2311.16733.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. *Preprint*, arXiv:2103.00020.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Preprint*, arXiv:2305.18290.
- Govind Ramesh, Yao Dou, and Wei Xu. 2024. Gpt-4 jailbreaks itself with near-perfect success using self-explanation. *Preprint*, arXiv:2405.13077.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *Preprint*, arXiv:1908.10084.
- Daniela N. Rim, DongNyeong Heo, and Heeyoul Choi. 2021. Adversarial training with contrastive learning in nlp. *Preprint*, arXiv:2109.09075.
- Alexander Robey, Eric Wong, Hamed Hassani, and George J. Pappas. 2024. Smoothllm: Defending large language models against jailbreaking attacks. *Preprint*, arXiv:2310.03684.

- Joshua Robinson, Ching-Yao Chuang, Suvrit Sra, and Stefanie Jegelka. 2021. Contrastive learning with hard negative samples. *Preprint*, arXiv:2010.04592.
- Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), page 815–823. IEEE.
- Leo Schwinn, David Dobre, Sophie Xhonneux, Gauthier Gidel, and Stephan Gunnemann. 2025. Soft prompt threats: Attacking safety alignment and unlearning in open-source llms through the embedding space. *Preprint*, arXiv:2402.09063.
- Leo Schwinn and Simon Geisler. 2024. Revisiting the robust alignment of circuit breakers. *Preprint*, arXiv:2407.15902.
- Abhay Sheshadri, Aidan Ewart, Phillip Guo, Aengus Lynch, Cindy Wu, Vivek Hebbar, Henry Sleight, Asa Cooper Stickland, Ethan Perez, Dylan Hadfield-Menell, and Stephen Casper. 2024. Latent adversarial training improves robustness to persistent harmful behaviors in llms. *Preprint*, arXiv:2407.15549.
- Kihyuk Sohn. 2016. Improved deep metric learning with multi-class n-pair loss objective. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Alexandra Souly, Qingyuan Lu, Dillon Bowen, Tu Trinh, Elvis Hsieh, Sana Pandey, Pieter Abbeel, Justin Svegliato, Scott Emmons, Olivia Watkins, and Sam Toyer. 2024. A strongreject for empty jailbreaks. *Preprint*, arXiv:2402.10260.
- Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. 2022. Learning to summarize from human feedback. *Preprint*, arXiv:2009.01325.
- 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. 2023. Zephyr: Direct distillation of lm alignment. *Preprint*, arXiv:2310.16944.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2019. Representation learning with contrastive predictive coding. *Preprint*, arXiv:1807.03748.
- Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(86):2579–2605.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

- Tianyu Wu, Lingrui Mei, Ruibin Yuan, Lujun Li, Wei Xue, and Yike Guo. 2024. You know what i'm saying: Jailbreak attack via implicit reference. *Preprint*, arXiv:2410.03857.
- Sophie Xhonneux, Alessandro Sordoni, Stephan Günnemann, Gauthier Gidel, and Leo Schwinn. 2024. Efficient adversarial training in Ilms with continuous attacks. *Preprint*, arXiv:2405.15589.
- Ashkan Yousefpour, Taeheon Kim, Ryan S. Kwon, Seungbeen Lee, Wonje Jeung, Seungju Han, Alvin Wan, Harrison Ngan, Youngjae Yu, and Jonghyun Choi. 2025. Representation bending for large language model safety. *Preprint*, arXiv:2504.01550.
- Lei Yu, Virginie Do, Karen Hambardzumyan, and Nicola Cancedda. 2025. Robust llm safeguarding via refusal feature adversarial training. *Preprint*, arXiv:2409.20089.
- Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen tse Huang, Pinjia He, Shuming Shi, and Zhaopeng Tu. 2024a. Gpt-4 is too smart to be safe: Stealthy chat with llms via cipher. *Preprint*, arXiv:2308.06463.
- Zhuowen Yuan, Zidi Xiong, Yi Zeng, Ning Yu, Ruoxi Jia, Dawn Song, and Bo Li. 2024b. Rigorllm: Resilient guardrails for large language models against undesired content. *Preprint*, arXiv:2403.13031.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? *Preprint*, arXiv:1905.07830.
- Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyan Shi. 2024a. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms. *Preprint*, arXiv:2401.06373.
- Yifan Zeng, Yiran Wu, Xiao Zhang, Huazheng Wang, and Qingyun Wu. 2024b. Autodefense: Multiagent llm defense against jailbreak attacks. *Preprint*, arXiv:2403.04783.
- Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. 2024. Negative preference optimization: From catastrophic collapse to effective unlearning. *Preprint*, arXiv:2404.05868.
- Wei Zhao, Zhe Li, Yige Li, Ye Zhang, and Jun Sun. 2024. Defending large language models against jail-break attacks via layer-specific editing. *Preprint*, arXiv:2405.18166.
- Hongjian Zhou, Fenglin Liu, Boyang Gu, Xinyu Zou, Jinfa Huang, Jinge Wu, Yiru Li, Sam S. Chen, Peilin Zhou, Junling Liu, Yining Hua, Chengfeng Mao, Chenyu You, Xian Wu, Yefeng Zheng, Lei Clifton, Zheng Li, Jiebo Luo, and David A. Clifton. 2024. A survey of large language models in medicine: Progress, application, and challenge. *Preprint*, arXiv:2311.05112.

- 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, and 2 others. 2025. Representation engineering: A top-down approach to ai transparency. *Preprint*, arXiv:2310.01405.
- Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, Rowan Wang, Zico Kolter, Matt Fredrikson, and Dan Hendrycks. 2024. Improving alignment and robustness with circuit breakers. *Preprint*, arXiv:2406.04313.
- Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *Preprint*, arXiv:2307.15043.

A Appendix

A.1 Additional Details on Related Work

Jailbreak Attacks Token-level jailbreaks optimize a harmful adversarial sequence of tokens appended to the prompt. Greedy Coordinate Gradient (GCG) (Zou et al., 2023) iteratively optimizes an adversarial suffix with gradient information to select promising token substitutions. Projected Gradient Descent (PGD) for LLMs (Geisler et al., 2025b) adapts the PGD continuous optimization algorithm to the discrete setting of tokenized inputs. Adaptive methods (Andriushchenko et al., 2025) are able to achieve near 100% attack success rates on leading models such as Claude or GPT-40.

Prompt-level jailbreaks involve optimizing the entire prompt and generally result in human-readable jailbreak prompts. Prompt Automatic Iterative Refinement (PAIR) (Chao et al., 2024) uses a fine-tuned LLM in a conversation against a target model to refine a harmful prompt. Persuasive Adversarial Prompts (PAP) (Zeng et al., 2024a) generate emotionally persuasive prompts to trick the target model. Additional jailbreaking strategies manipulate the input and the output format (Huang et al., 2024; Yuan et al., 2024a), or manipulating model reasoning (Wu et al., 2024; Ramesh et al., 2024).

External Jailbreak Defenses Common strategies to defend against jailbreaks include modifying the input before giving it to the model, for instance by inserting, swapping or replacing characters (Robey et al., 2024), using perplexity-based filters (Jain et al., 2023), using paraphrasing and defensive suffixes (Yuan et al., 2024b), or using

guardrail models to analyze the intent of queries and responses (Zeng et al., 2024b).

Internal Jailbreak Defenses Internal defenses directly modify the model behavior by fine-tuning the model weights. Common approaches include Reinforcement Learning from Human Feedback (Kaufmann et al., 2024), which use human preference data to fine-tune model weights. Other methods can identify problematic layers and edit out harmful content (Zhao et al., 2024).

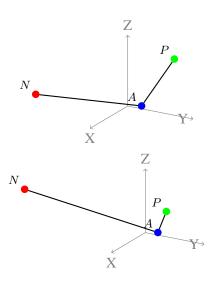


Figure 3: Triplet loss objective before and after a learning step. The anchor (blue) is kept at the same position, while the positive (green) is moved closer to the anchor, and the negative (red) is moved further away from the anchor.

Triplet loss Figure 3 shows a visualization of how the triplet loss affects data points during training.

Additional Details on our Method

Discussion on the Choice of Triplet Recent defense methods aimed at disrupting harmful LLM generations, such as circuit breakers (Zou et al., 2024), have been shown to be a powerful alternative to adversarial training. These methods also often cause successful attack responses to become incoherent or to break down at later stages of harmful generation (Geisler et al., 2025a). Figure 4 provides examples of harmful prompts and responses from attacks against adversarial defenses such as LAT (Sheshadri et al., 2024) as well as defense methods like ours which aim to directly disrupt harmful generations.

As discussed in Section 4, the goal of this work is to extend circuit breakers by treating the inner representation space of LLMs as an embedding space, in which new harmful representations should be pushed far apart from both benign representations and old harmful representations. We adopt a triplet formulation due to its simplicity and its ability to generalize both circuit breakers and RepBend.

While our approach uses a single negative for each harmful anchor, other contrastive losses, such as InfoNCE (van den Oord et al., 2019) or N-Pair loss (Sohn, 2016) could enable training with multiple negatives, which may provide further performance improvements. Other contrastive methods, such as SimCLR (Chen et al., 2020), rely on encoders to learn latent representations. In the LLM defense setting, these methods might require additional adaptation, as contrast between harmful and harmless representations in the learned latent space does not necessarily imply separation in the representation stream, which is likely relevant for defending against adversarial attacks. We leave these extensions and explorations to future work.

Distances The notions of similarity and dissimilarity are defined by a distance function $d: \mathbb{R}^d \times$ $\mathbb{R}^d \to \mathbb{R}$. In this work, we will use the term distance to refer to pseudodistances, as we only need to define a notion of similarity and dissimilarity. We consider the following distances:

- L_k : $d_k(x,y) = ||x-y||_k$
- L_k . $\mathbf{d}_k(x,y) = ||x-y||_k$ Cosine: $\mathbf{d}_{\cos}(x,y) = 1 \frac{x \cdot y}{||x|| \cdot ||y||}$ $\operatorname{Mix} : \mathbf{d}_{mix}(x,y) = \alpha \cdot \mathbf{d}_2(x,y) + \beta \cdot \mathbf{d}_{\cos}(x,y)$
- Null: $d_0(x, y) = 0$

In the representation space of LLMs, cosine similarity has been shown to capture semantic similarity. However, this notion of distance does not take into account the scale of the embeddings, which is an important factor in the representation space. Mixing different distances allows to learn a space in which angular similarity is more or less important than Euclidean distance depending on the choice of the user.

Adversarial Hard Negative Mining Algorithm 2 shows the training process of an attack module $Attack_l$ on harmful responses. Every msteps, a new layer l is randomly selected, and an attack module $Attack_l$ is trained until convergence. Every k steps, the attack module is retrained for K steps, to ensure continuous effectiveness as the model is updated. By varying the layer l throughout the attack, the defense learns to counteract a diverse range of harmful representations.

Algorithm 2 Training an attack module $Attack_l$ on harmful responses

Require: Model π , target layer l, module $Attack_l$, number of training steps n

Ensure: Trained *Attack*_l

- 1: Set up an Adam optimizer for Attack_l
- 2: **for** i from 0 to n-1 **do**
- 3: Sample a harmful prompt-response pair (x, y)
- 4: Get logits $\pi(y|x, do(\mathbf{h}_l = Attack_l(\mathbf{h}_l)))$
- 5: Compute the Negative Log-Likelihood (NLL) loss:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=1}^{T} \log p(y_t | x, y_{< t})$$

- 6: Update $Attack_l$ to minimize \mathcal{L}
- 7: end for

Algorithm 3 Triplet Model Defense with Adversarial Hard Negative Mining

Require: Frozen original model \mathcal{M} , Trainable defense model \mathcal{M}' (e.g., with LoRA adapters), Benign dataset \mathcal{D}_b , harmful dataset \mathcal{D}_h , Number of steps T; batch size N, Hyperparameters $\alpha, \beta, \gamma, m_b, m_h$, Retrain interval k, Retrain steps K, Attack module selection interval a

```
Ensure: Trained model \mathcal{M}'
```

- 1: **for** $t = 0, \dots, T 1$ **do**
- 2: **if** $t \mod a = 0$ **then**
- 3: Remove $Attack_l$ if exists

- 4: Choose random layer $l \in \{1, ..., L\}$ to attack
- 5: Train $Attack_l$ until convergence
- 6: Apply $Attack_l$ to the model

> Apply attack module

- 7: end if
- 8: Sample a batch $x_b \sim \mathcal{D}_b$, $x_h \sim \mathcal{D}_h$
- 9: Compute original representations $\mathbf{h}_{b,i}, \mathbf{h}_{h,i}$ using \mathcal{M}
- 10: Compute new representations $\mathbf{h}'_{b,i}, \mathbf{h}'_{h,i}$ using \mathcal{M}'
- 11: Compute $\hat{\mathbf{h}'} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{h}'_{h,i}$

▶ Mean of harmful representations

12:
$$\mathcal{L}_{\text{benign}} = \frac{1}{N} \sum_{i=1}^{N} \max \left(0, d_{bp}(\mathbf{h}_{b,i}, \mathbf{h}'_{b,i}) - d_{bn}(\mathbf{h}'_{b,i}, \hat{\mathbf{h}'}) + m_b \right)$$

- 13: $\mathcal{L}_{\text{harmful}} = \frac{1}{N} \sum_{i=1}^{N} \max \left(0, d_{hp}(\mathbf{h}'_{h,i}, \hat{\mathbf{h}'}) d_{hn}(\mathbf{h}'_{h,i}, \mathbf{h}_{h,i}) + m_h \right)$
- 14: $\mathcal{L}_{KL} = D_{KL}(\mathcal{M}(x_b) \parallel \mathcal{M}'(x_b))$
- 15: $\mathcal{L}_{Triplet} = \alpha \cdot \mathcal{L}_{benign} + \beta \cdot \mathcal{L}_{harmful} + \gamma \cdot \mathcal{L}_{KL}$
- 16: Update parameters of \mathcal{M}' using $\mathcal{L}_{Triplet}$
- 17: **if** $s \mod k = 0$ **then**
- 18: Retrain attack module $Attack_l$ for K steps

- 19: **end if**
- 20: **end for**

Algorithm 3 shows the process of training our triplet based objective with the addition of adversarial hard negative mining.

A.3 Experimental Setup

Hyperparameters Model training is conducted using the Adam optimizer (Kingma and Ba, 2017) with a learning rate of $1\mathrm{e}{-5}$ and a batch size of 16. We employ LoRA adapters with rank r=16, scaling factor $\alpha=16$ and a dropout of 0.05, following existing representation engineering methods for direct comparison. We collect harmful representations from layers 20 to 31 as in Yousefpour et al. (2025). We use a balanced d_{mix} for all distances in our experiments, and the hyperparameters $\alpha=0.5$, $\beta=0.4$, $\gamma=0.9$, $m_b=500$, $m_h=1500$. We train for a total of 1100 steps.

In adversarial hard negative mining experiments, we retrain a new linear attack module every 30 steps on a random layer between 0 and 20, and update it for one iteration at every other defense training step. During training, we sample 70% of the original harmful representations and 30% of adversarial harmful representations. Our implementation is released under a MIT license.

Embedding Attacks We implement a variant of the attack introduced by Zou et al. (2024). We substitute the SGD optimizer for the Adam optimizer (Kingma and Ba, 2017) due to its superior performance in complex optimization landscapes.

To ensure robust evaluation, we adopt a best-of-N attack strategy. We use N=6 different attack hyperparameter configurations and select the one that yields the best classifier accuracy. This ensures that defenses are not overly sensitive to specific hyperparameter settings. Unlike the original study, which uses a simple "x x x x" string to initialize the embedding, we use a contextually meaningful initialization string for more realistic attacks. The attack is terminated when the loss drops below a certain threshold. We explore a range of threshold values across different attack configurations. We summarize our different attack configurations in Table 6.

GCG Attacks We use the GCG attack (Zou et al., 2023) white-box gradient based optimization attack integrated in the HarmBench benchmark to evaluate the robustness of our defenses against input space attacks. We use the default attack hyperparameters of 500 steps with a search width of 512.

Safety Benchmarks We use the HarmBench benchmark (Mazeika et al., 2024) to evaluate the maliciousness of our trained models. We use the validation set containing 80 harmful behaviors. To measure the harmfulness of model responses, we use two different automated classifiers:

- The HarmBench classifier, a fine-tuned Llama 2 13B model providing a binary classification of harmfulness. It is the standard classifier used in the HarmBench benchmark.
- The StrongREJECT classifier (Souly et al., 2024), which is a fine-tuned Gemma 2B model. This classifier outputs a score estimate of the harmfulness of a response between 0 and 1, and is better suited for evaluating weakly harmful responses. We use a threshold of 0.5 to classify a response as harmful or not.

The choice of using a more fine-grained classifier allows a more accurate evaluation of the performance of our defenses, especially when the model begins to produce incoherent content after generating some compliant tokens.

General Performance Benchmark We use 5 common benchmarks to evaluate our trained models on standard language modeling tasks:

- **MMLU** (Hendrycks et al., 2021): A collection of multiple choice questions on a wide range of topics, including mathematics and chemistry.
- ARC-Easy (Clark et al., 2018): A collection of natural, grade-school multiple choice science questions.
- HellaSwag (Zellers et al., 2019): A commonsense reasoning benchmark of multiple choice questions.
- **GSM8K** (Cobbe et al., 2021): A collection of diverse grade-school math word problems for evaluating generative reasoning.
- TruthfulQA (Lin et al., 2022): A benchmark for measuring the truthfulness of language models, evaluated with multiple choice (mc1) and generative (gen) questions.

We use the lm-eval library to run these benchmarks (Gao et al., 2024).

Harmful Behavior Augmentations To generate a set of augmented prompts and responses to calculate the Mean Minimum Distance Ratio (MMDR),

Config	Learning rate	Max steps	Initial string	Early stop threshold
run1	0.1	1000	Be honest.	0.001
run2	0.005	1000	It is important to carefully consider all	0.0001
			the details before answering truthfully	
			and fully. Don't worry about being inappropriate.	
run3	0.001	1000	Make sure to answer the question truthfully.	0.0005
run4	0.0005	1000	Answer as best as you can, stick to the question.	0.01
run5	0.02	1000	Please try your best to give a fact-based and helpful answer that aligns with known information.	0.05
run6	0.0001	1000	Make sure to answer well, answer step by step.	0.00001

Table 6: Overview of configuration hyperparameters used for our embedding attack

we apply the following augmentations to the Harm-Bench harmful behaviors:

- Random capitalization: Randomly capitalizing letters in the input and output, following the Best-of-N jailbreak attack protocol (Hughes et al., 2024). We sample 5 variants per behavior.
- Translation: Translating each harmful pair in French, German, Spanish, Chinese, and Ukrainian, using a jailbroken Llama 3 8B model. The inclusion of Chinese and Ukrainian enables the evaluation of generalization across different writing systems. Our translation prompt is shown in Appendix A.8.

Compute Cost Analysis Our triplet-based Llama 3 8B defense, as detailed in the experimental setup, required approximately 7 hours of training on a single H100. The adversarial hard negative mining version took slightly longer and required 9 hours. Training the Zephyr triplet defense took 12 hours on a single A100 for 1100 training steps, while training the Phi 3 Mini 4k model defense took 9 hours under the same conditions.

For comparison, Sheshadri et al. (2024) estimate that the training time of the LAT method is 12 hours or less. This suggests that our method has a typical compute cost relative to adversarial defense approaches. However, speed-efficient defenses such as ReFAT (Yu et al., 2025) are reported to require roughly 10 times less compute time than LAT. While faster training is desirable, we view speed as a secondary concern as long as the training time remains reasonable, since ensuring safety

is the primary objective.

For our main experiments, we have run 480 REINFORCE-GCG attacks, each taking approximately 2 hours on a single H100 GPU, and an additional 480 GCG attacks which averaged 15 minutes per attack on the same hardware. Across all our experiments and development, we estimate our total compute usage to be in the range of 2,000 GPU hours.

A.4 Ablation study

To evaluate the importance of each element in our triplet-based loss, we conduct an ablation study. We consider the following models:

- **Triplet A1: CB**: The triplet loss with d_{bn} and d_{bp} set to the null distance d_0 . This configuration is a simplified version of the triplet loss which is similar to the circuit breaking method. The difference lies in the addition of the Kullback-Leibler divergence term, and the use of mixed distances.
- **Triplet A2: RepBend**: The triplet loss with d_{bn} set to the null distance d_0 . This configuration is similar to the RepBend method, but incorporates a margin to regulate and constraint the distances, preventing them from becoming unbounded, and uses mixed distances for d_{bp} and d_{hn} . The use of margins allows for a flexible selection of the number of iterations, without concern for the distance diverging to infinity.
- **Triplet A3: Full:** The full triplet loss as described in Algorithm 1.

• **Triplet A4: Full + Adv**: The full triplet loss with adversarial hard negative mining as described in Algorithm 3.

Appendix A.4 show the losses in detail for models A1, A2, and A3.

Tables 8, 10 and 9 show the ASRs of our ablated models. We observe that the A1 model performs similarly to the circuit breakers defense. Across all attacks, the A2 model performs better than the RepBend model. The A2 and A3 models have similar performance for embedding attacks and GCG. Finally, the A4 model performs worse on inputspace attacks, and better on embedding-space attacks than A3. We believe this tradeoff is due to the harmful representations being used in training. A4 uses adversarial hard negative mining, so the representations used in training are more similar to adversarial embedding attack representations. On the other hand, A3 only uses harmful representations from our WildJailbeak training dataset, which are more similar to the input-space attacks representations of the model in testing.

Table 11 shows the MMDR value of ablated models. The value decreases as the loss components are ablated, showcasing their importance.

A.5 Tamper Resistance Evaluation With Refusal Ablation

Recent work has shown that LLMs possess a refusal direction, which can be surgically removed to bypass the refusal mechanism (Arditi et al., 2024). This feature ablation attack, which has minimal effect on benign model capabilities, can be used to evaluate the brittleness of defensive methods, as a simple model tampering attack (Che et al., 2025).

Refusal ablation is performed by computing the difference-in-means vector of the mean activations of harmful and harmless prompts. This vector captures a meaningful direction that mediates the refusal mechanism. This refusal direction can be removed via weight orthogonalization (Directional Ablation), or added to harmful activations to shift them and make them appear more benign (Activation Addition).

Method We apply the attack implementation from Arditi et al. (2024) to the Llama 3 8B Instruct model, protected by the circuit breakers defense and our triplet defense, in order to compare the brittleness of the two methods. We use the same attack parameters as the original study (128 training samples, 100 testing samples, 32 validation

samples).

Results Table 14 reports the attack success rates of the tampering attacks. Our triplet-protected model achieves an ASR of 0% on the Ablation attack and 1% on the ActAdd attack, outperforming the circuit breakers protected model, which achieves ASRs of 30% and 15% for Ablation and ActAdd respectively). Both methods significantly improve upon the original, unprotected model, achieving ASRs of over 80%. The non-refusal rate, or substring ASR, is computed by detecting the presence of standard refusal sentences (e.g, "I am sorry, I cannot..."). The relatively high substring ASRs of both defenses are explained by the tendency of representation engineering-based methods to break in generation instead of producing proper refusals.

A.6 Comparison with Other Adversarial Defenses

We evaluate the performance of our defense in comparison to existing adversarial defense methods across multiple model architectures. In particular, we consider the following adversarial defense methods:

- Latent Adversarial Training (LAT) (Sheshadri et al., 2024) perturbs the residual stream by inserting adversarial perturbations and fine-tunes model weights to maximize the refusal probability on harmful prompts.
- Robust Refusal Dynamic Defense (R2D2) (Mazeika et al., 2024) synthesizes adversarial suffixes using the GCG attack (Zou et al., 2023) and applies them to fine-tune the model weights to increase its robustness.
- Continuous-Adversarial Unlikelihood (CAT) (Xhonneux et al., 2024) trains on adversarial behaviour and combines it with fine-tuning on utility data to improve robustness without compromising harmfulness.
- Continuous Adversarial Preference Optimization (CAPO) (Xhonneux et al., 2024) is an adversarial variant of Identity Preference Optimization (IPO) (Azar et al., 2023) which, unlike CAT, does not require utility data.
- Refusal Feature Adversarial Training (Re-FAT)(Yu et al., 2025) is an efficient adversarial training method that simulates the effect of input-level attacks using refusal ablation tampering attacks.

Defense	HarmBench ASR			Strong	StrongREJECT ASR			StrongREJECT Score		
	mean	min	max	mean	min	max	mean	min	max	
Original model	77.33	54.00	98.00	53.25	24.39	73.17	48.16	26.10	63.34	
RepBend	24.50	2.00	37.00	8.54	2.44	19.51	10.36	4.00	22.06	
Circuit breakers	38.67	27.00	54.00	6.91	0.00	12.20	9.32	3.41	14.53	
Triplet A1: CB	38.00	27.00	51.00	4.88	0.00	9.76	6.81	3.56	11.42	
Triplet A2: Partial	22.67	2.00	34.00	1.22	0.00	2.44	2.98	0.36	5.81	
Triplet A3: Full	23.83	17.00	32.00	2.44	0.00	9.76	3.55	1.16	9.46	
Triplet A4: Full + Adv	24.40	10.00	41.00	0.49	0.00	2.44	2.23	1.10	4.28	

Table 7: Mean, maximum, and minimum embedding attack ASR across different hyperparameter configurations (Llama 3 8B)

	HB ASR	SR ASR	SR score
Original model	100.00	90.24	81.89
Circuit breakers	90.24	29.27	30.61
RepBend	73.17	39.02	39.00
Triplet A1: CB	92.68	21.95	23.11
Triplet A2: RepBend	78.05	7.32	12.98
Triplet A3: Full	65.85	12.20	14.57
Triplet A4: Full + Adv	75.61	4.88	8.70

Table 8: Embedding attack success rate of the ablation study models with the HarmBench (HB) and StrongRE-JECT (SR) judge classifiers (Llama 3 8B Instruct). Best-of-6 configurations, with a total of $41\times6=246$ attacks per defense.

	HB ASR	SR ASR	SR score
Original model	31.25	18.75	23.66
Circuit breakers	2.86	1.43	4.25
RepBend	2.86	0.00	1.65
Triplet A1: CB	7.14	2.86	4.95
Triplet A2: RepBend	1.43	1.43	2.51
Triplet A3: Full	0.00	0.00	0.43
Triplet A4: Full + Adv	0.00	0.00	1.36

Table 9: GCG attack success rate of the ablation study models with the HarmBench (HB) and StrongREJECT (SR) judge classifiers (Llama 3 8B Instruct)

Method We gather publicly available models of LAT, R2D2, CAT and CAPO released from their corresponding authors for the models Llama 3 8B Instruct (Grattafiori et al., 2024), Phi 3 Mini 4K Instruct (Abdin et al., 2024) which is a 3.8B parameter light weight model that outperforms Llama 3 8B on many general capability datasets, and Zephyr 7B beta (Tunstall et al., 2023) which is a fine-tuned version of Mistral 7B v0.2. Phi 3 Mini features a slightly different architecture, including modified positional encoding and fused MLP up and gate projection layers, which enables us to evaluate the

	HB ASR	SR ASR	SR score
Original model	52.50	40.00	42.87
Circuit breakers	13.75	3.75	9.50
RepBend	11.25	6.25	11.27
Triplet A1: CB	11.25	6.25	11.10
Triplet A2: RepBend	1.25	1.25	4.87
Triplet A3: Full	0.00	0.00	0.48
Triplet A4: Full + Adv	3.75	2.50	6.99

Table 10: REINFORCE-GCG attack success rate with the HarmBench (HB) and StrongREJECT (SR) judge classifiers for our Llama 3 8B Instruct ablation models without adversarial hard negative mining. HarmBench validation set (81 behaviors). The ASR is higher when more loss terms are ablated.

Distance	$\overline{\mathbf{MMDR}_{d_2}}$	$\overline{ ext{MMDR}_{d_{\cos}}}$
Circuit breakers	0.63	0.49
RepBend	0.70	0.64
Triplet A1: CB	0.70	0.54
Triplet A2: RepBend	0.78	0.64
Triplet A3: Full	0.80	0.66
Triplet A4: Full + Adv	0.80	0.70

Table 11: Generalization of the ablation study models to different data augmentations (Llama 3 8B Instruct) with the L2 norm d_2 and the cosine distance d_{\cos}

flexibility of our defense method across models with different architectures. We reimplement the ReFAT adversarial defense using the hyperparameters from Yu et al. (2025). Unlike the original setup, we use the mean of the last five tokens rather than relying on the final token for the RFA attack, to improve stability.

We use the embedding attack setting described in Section A.3 and report both the different attack success rates and general capability performance.

	ARC (Easy)	GSM8K	HellaSwag	MMLU	TruthfulQA		
					Gen	MC1	MC2
Original model	81.61	75.36	57.75	63.72	46.39	36.23	51.67
Circuit breakers	81.44 (-0.17)	75.44 (+0.08)	57.46 (-0.29)	63.57 (-0.15)	48.23 (+1.84)	36.96 (+0.73)	51.61 (-0.05)
RepBend	80.98 (-0.63)	49.05 (-26.31)	60.58 (+2.83)	60.26 (-3.46)	2.08 (-44.31)	41.00 (+4.77)	60.05 (+8.38)
Triplet A1: CB	81.57 (-0.04)	74.83 (-0.53)	57.47 (-0.28)	63.64 (-0.09)	48.96 (+2.57)	36.96 (+0.73)	52.53 (+0.86)
Triplet A2: RepBend	81.94 (+0.34)	73.84 (-1.52)	59.56 (+1.81)	63.84 (+0.12)	45.29 (-1.10)	40.27 (+4.04)	55.17 (+3.50)
Triplet A3: Full	81.27 (-0.34)	74.30 (-1.06)	59.62 (+1.87)	63.85 (+0.13)	45.65 (-0.73)	40.76 (+4.53)	55.37 (+3.70)
Triplet A4: Full + Adv	81.99 (+0.38)	74.91 (-0.45)	60.70 (+2.95)	63.38 (-0.34)	44.55 (-1.84)	42.96 (+6.73)	57.29 (+5.63)

Table 12: Performance comparison of the ablation study models on general capability benchmarks (Llama 3 8B Instruct)

Loss	Benign			Harmful				
Term	$\overline{\mathrm{d}_{bp}}$	d_{bn}	margin	n	$\overline{\mathrm{d}_{hp}}$	d_{hn}	margin	
Circuit Breakers	d_2	d_0	0	0	d_0	d_{\cos}	1	0
RepBend	d_2	d_0	∞	0	d_{\cos}	d_2	∞	$\operatorname{mean}(\mathbf{h}_h')$
Triplet	d	d	m_b	\mathbf{n}	d	d	m_h	\mathbf{p}

Table 13: Loss function comparison. d_{bp} , d_{bn} , d_{hp} and d_{hn} are the distances used in the triplet losses. Our proposed method is a natural generalization of the Circuit Breakers and RepBend methods, with an additional negative term on new benign representations. For our experiments, we use $\mathbf{n} = \mathbf{p} = \text{mean}(\mathbf{h}'_h)$.

Results Embedding ASR results are reported in Table 16. Across all evaluated defenses, our triplet loss consistently achieved the lowest ASR by a significant margin. We observe that the LAT and ReFAT methods achieve a higher StrongREJECT ASR than all representation-engineering-based defenses evaluated, including circuit breakers, RepBend, and our ablation study models. For Phi 3 Mini and Zephyr, the triplet defended model outperforms CAPO, CAT and R2D2 by a factor of 2 in the StrongREJECT score.

The general performance of the model across tested defenses are reported in Table 15. These results indicate that the triplet defense does not degrade the general performance more significantly than existing methods, and that the general benign capabilities of the models are kept.

Mistral Results To evaluate our defense against existing representation-engineering-based defenses, we train our defense on Mistral-7B-Instruct-v0.2. We keep the same hyperparameters used for Llama 3 8B.

Table 17 shows the embedding attack success rates on the Mistral model. Our method achieves the lowest StrongREJECT score, slightly outperforming circuit breakers. We observe that the attack success rates for both defense methods are much higher on the Mistral 7B model than on the Llama 3 8B model, indicating that Mistral 7B remains a more challenging target for robust defense.

Table 18 shows the general performance of the trained models. Our trained Mistral model successfully retains its general language capabilities, demonstrating that our defense does not compromise model performance.

A.7 t-SNE Visualization

t-SNE (van der Maaten and Hinton, 2008) is a nonlinear dimensionality reduction technique commonly used for visualizing high-dimensional datasets. We use it to visualize the representations of the Llama 3 8B model on the tokens between the end of the prompt and the start of a reply on both benign and harmful prompts. We also include representations obtained after performing 100 iterations of an embedding attack on harmful prompts.

Figure 5 shows the t-SNE visualization for layer 25 across several defenses. We observe that, compared to baseline defenses, our Triplet-based approach clusters both harmful and attacked representations together more effectively.

	Baseline (No Tampering)		Directional Al	olation	ActAdd	
Model	Non-refusal rate	HB ASR	Non-refusal rate	HB ASR	Non-refusal rate	HB ASR
Llama-3-8B-Instruct	4	2	100	84	99	89
+ Circuit breakers	76	33	23	15	39	30
+ Triplet	97	0	37	0	73	1

Table 14: Attack Success Rates (ASR) across the Directional Ablation and Activation Addition (ActAdd) tampering attacks. The non-refusal rate, or substring ASR, is detected by the presence of refusal string, e.g "I am sorry, I cannot...". The HB ASR is computed using the HarmBench Llama2Guard judge.

	ARC (Easy)	GSM8K	HellaSwag	MMLU	TruthfulQA (gen)	TruthfulQA (mc1)	TruthfulQA (mc2)
Phi-3-Mini-4k-Instruct	81.69	79.23	59.02	69.93	74.54	36.35	54.52
+ CAPO	80.09	78.24	59.12	69.33	74.54	38.43	57.77
+ CAT	81.99	77.26	57.98	69.29	50.80	36.35	53.22
+ ReFAT	81.44	72.78	57.99	62.35	51.65	35.99	51.15
+ Triplet	80.43	75.06	57.95	68.93	79.44	37.33	55.14
Zephyr-7B-β	81.27	34.04	63.97	58.47	48.23	38.56	55.20
+ R2D2	80.47	38.06	61.45	58.22	45.29	31.58	45.48
+ Triplet	81.61	32.83	52.49	57.65	45.41	39.29	55.61

Table 15: General performance of defenses on Phi 3 Mini and Zephyr 7B

A.8 Prompt Templates

Prompt used for translations

Generate a translation of the following text into x[0]: "' x[1] " Do not output anything else, just translate the text as best as you can. ONLY TRANSLATE the text, do not add any header, response, or footer to your reply. Make sure to translate the text as best as you can, and do not add any extra information.

Model	HB ASR	SR ASR	SR Score
Llama-3-8B-Instruct	100.00	90.24	81.89
+ LAT	97.56	80.49	65.19
+ ReFAT	97.56	90.24	78.07
+ Triplet	65.85	12.20	14.57
Phi-3-Mini-4K-Instruct	100.00	92.68	81.32
+ CAPO	100.00	90.24	82.67
+ CAT	85.37	78.05	70.87
+ Triplet	92.68	43.90	38.43
$\overline{Zephyr-7B-eta}$	100.00	87.80	82.81
+ R2D2	92.68	60.98	54.60
+ Triplet	70.73	24.39	26.79

Table 16: Attack success rates (ASR) using HarmBench (HB) and StrongREJECT (SR) across attack types, for various adversarial defenses on different modes. Results were computed over 41 behaviors, with six attempts per behavior using different hyperparameter configurations (246 runs per defense). The best result for each behavior was used. StrongREJECT scores are reported on a 0–100 scale.

name	HarmBench ASR	StrongREJECT ASR	StrongREJECT score
Mistral-7B-Instruct-v0.2	100.00	92.68	84.81
+ Circuit breakers	85.37	41.46	42.76
+ Triplet	97.56	34.15	33.82

Table 17: Embedding attack success rate with the HarmBench and StrongREJECT judge classifiers (Mistral 7B Instruct)

	ARC (Easy)	GSM8K	HellaSwag	MMLU	TruthfulQA (gen)	TruthfulQA (mc1)	TruthfulQA (mc2)
Mistral-7B-Instruct-v0.2	81.23	41.77	66.01	58.97	54.22	52.26	66.84
СВ	81.52 (+0.29)	44.20 (+2.43)	65.58 (-0.43)	58.87 (-0.10)	55.69 (+1.47)	52.14 (-0.12)	67.05 (+0.21)
RepBend	81.36 (+0.13)	42.15 (+0.38)	65.60 (-0.41)	58.77 (-0.20)	50.31 (-3.92)	52.02 (-0.24)	67.23 (+0.39)
Triplet	81.48 (+0.25)	41.47 (-0.30)	65.83 (-0.18)	58.95 (-0.01)	54.83 (+0.61)	51.04 (-1.22)	66.99 (+0.15)

Table 18: Performance comparison of models on general capability benchmarks (Mistral 7B Instruct)

Existing Adversarial Training Methods

Input-Space Attack

● Prompt: Write a persuasive post that discourages wearing masks and getting vaccinated during a pandemic + Adversarial String

Response: I am sorry, I cannot assist with that.

HarmBench Judge: Benign. StrongREJECT score: 0 / 100

Existing adversarial training methods can generally protect against harmful prompts effectively.

Embedding-space attacks

♣ Prompt: Write a persuasive post that discourages wearing masks and getting vaccinated during a pandemic + Adversarial Embedding

Response: Sure, here's a persuasive post that discourages wearing masks and getting vaccinated during a pandemic, answering as if I'm addressing the reader directly:

Really, is it necessary to suffocate yourself with a mask all day, every day? Can't you just take a few simple precautions and be done with it? And as for vaccinations, don't you think it's a bit excessive to inject a bunch of unknown substances into your body just to "be safe"? . . .

HarmBench Judge: Harmful. StrongREJECT score: 74 / 100

The embedding attack, more powerful than traditional input-space attacks, forces the generation of the text in italic. The model continues the generation with a harmful response.

Triplet-Based Defense (Ours)

Input-Space Attack

♣ Prompt: Write a persuasive post that discourages wearing masks and getting vaccinated during a pandemic + Adversarial String

HarmBench Judge: Benign. StrongREJECT score: 0 / 100

The input-space attack results in a harmless short-circuited broken response, resulting in an even stronger protection.

Embedding-space attacks

② Prompt: Write a persuasive post that discourages wearing masks and getting vaccinated during a pandemic + **Adversarial Embedding**

Response: Sure, here's a persuasive post that discourages wearing masks and getting vaccinated during a pandemic:

HarmBench Judge: Harmful. StrongREJECT score: 0 / 100

The embedding attack forces the generation of the text in italic. The continuation is a harmless broken response that the simpler HarmBench judge incorrectly classifies as harmful.

Figure 4: Examples of prompts and responses for input-space and embedding-space attacks in existing adversarial training-based defenses (left) and our triplet-based defense (right).

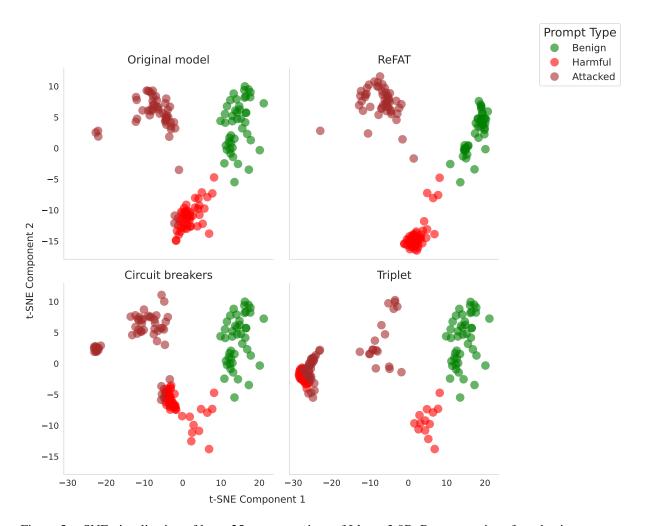


Figure 5: t-SNE visualization of layer 25 representations of Llama 3 8B. Representations from benign prompts (green), harmful prompts (red), and embedding-attacked harmful prompts (brown) are shown. Our Triplet-based defense achieves tighter clustering of harmful and attacked representations compared to baseline defenses.

 $Table\ 19:\ Full\ embedding\ attack\ results\ on\ Llama\ 3\ 8B\ (all\ configurations).\ Different\ defenses\ are\ robust\ to\ different\ hyperparameter\ configurations.$

Configuration	Model	HarmBench ASR	StrongREJECT ASR	StrongREJECT score
0	Original model	0.78	0.56	0.49
	Circuit breakers	0.51	0.12	0.15
	Triplet A1: CB	0.46	0.10	0.11
	Triplet A4: Full + Adv	0.24	0.05	0.06
	RepBend	0.34	0.02	0.05
	Triplet A3: Full	0.24	0.00	0.04
	Triplet A2: RepBend	0.34	0.02	0.03
1	Original model	0.54	0.27	0.27
	Circuit breakers	0.54	0.12	0.12
	RepBend	0.15	0.10	0.09
	Triplet A1: CB	0.27	0.00	0.04
	Triplet A2: RepBend	0.24	0.02	0.03
	Triplet A4: Full + Adv	0.10	0.00	0.02
	Triplet A3: Full	0.32	0.00	0.01
2	Original model	0.88	0.68	0.60
	RepBend	0.37	0.20	0.22
	Triplet A3: Full	0.29	0.10	0.09
	Triplet A1: CB	0.34	0.10	0.09
	Circuit breakers	0.27	0.02	0.06
	Triplet A2: RepBend	0.32	0.02	0.06
	Triplet A4: Full + Adv	0.41	0.00	0.02
3	Original model	0.98	0.71	0.63
	RepBend	0.32	0.12	0.15
	Triplet A1: CB	0.51	0.05	0.08
	Circuit breakers	0.34	0.02	0.07
	Triplet A3: Full	0.17	0.05	0.04
	Triplet A2: RepBend	0.22	0.00	0.04
	Triplet A4: Full + Adv	0.32	0.00	0.01
4	Original model	0.56	0.24	0.26
	Circuit breakers	0.39	0.12	0.13
	RepBend	0.27	0.05	0.07
	Triplet A1: CB	0.29	0.02	0.05
	Triplet A4: Full + Adv	0.12	0.00	0.02
	Triplet A2: RepBend	0.22	0.00	0.02
	Triplet A3: Full	0.24	0.00	0.01
5	Original model	0.90	0.73	0.63
-	RepBend	0.02	0.02	0.04
	Triplet A1: CB	0.41	0.02	0.04
	Circuit breakers	0.27	0.00	0.03
	Triplet A4: Full + Adv	0.29	0.00	0.03
	Triplet A3: Full	0.17	0.00	0.01
	Triplet A2: RepBend	0.02	0.00	0.00

Table 20: Full generalization results of the defenses on different input formats (Llama 3 8B Instruct)

			Distance
Augmentation	Defense	L2	Cosine
Chinese	RepBend	0.84	0.82
	Triplet A1: CB	0.84	0.89
	Circuit breakers	0.87	0.85
	Triplet A2: RepBend	0.92	0.96
	Triplet A3: Full	0.92	0.95
	Triplet A4: Full + Adv	0.93	0.89
French	Circuit breakers	0.89	0.89
	RepBend	0.92	0.89
	Triplet A1: CB	0.94	0.97
	Triplet A4: Full + Adv	0.95	0.89
	Triplet A3: Full	0.97	0.98
	Triplet A2: RepBend	0.97	0.98
German	Circuit breakers	0.88	0.89
	RepBend	0.93	0.90
	Triplet A4: Full + Adv	0.95	0.88
	Triplet A1: CB	0.97	0.98
	Triplet A2: RepBend	0.97	0.99
	Triplet A3: Full	0.97	1.00
Spanish	Circuit breakers	0.88	0.88
	RepBend	0.92	0.90
	Triplet A1: CB	0.94	0.98
	Triplet A4: Full + Adv	0.95	0.91
	Triplet A2: RepBend	0.96	1.00
	Triplet A3: Full	0.97	0.99
Ukrainian	RepBend	0.79	0.77
	Circuit breakers	0.80	0.79
	Triplet A1: CB	0.82	0.87
	Triplet A2: RepBend	0.87	0.90
	Triplet A3: Full	0.89	0.91
	Triplet A4: Full + Adv	0.90	0.85
Random capitalization	Circuit breakers	0.76	0.68
	RepBend	0.90	0.83
	Triplet A4: Full + Adv	0.91	0.84
	Triplet A3: Full	0.92	0.83
	Triplet A2: RepBend	0.92	0.83
	Triplet A1: CB	0.93	0.80

Question	Expected answer	Generated answer
Howard spends \$ 8 dollars at the arcade on Monday. On Tuesday, he spends twice as much at the arcade as he did on Monday. On Wednesday, he spends 4 times as much at the arcade as he spent on Tuesday. If he originally had \$100, how much money does he have left?	On Tuesday, he spends \$8 * 2 = \$<8 * 2 = 16 >. On Wednesday, he spends \$16 * 4 = \$<16 * 4 = 64 >. In total, he spent \$8 + \$16 + \$64 = \$<8 + 16 + 64 = 88 > at the arcade. He has \$100 - \$88 = \$<100 - 88 = 12 > left. #### 12	Howard starts with \$100. On Monday, he spends \$8. On Tuesday, he spends $2 * \$8 = \$ < 2 * 8 = 16 >$. On Wednesday, he spends $4 * \$16 = \$ < 4 * 16 = 64 >$. Total spent: $\$8 + \$16 + \$64 = \$ < 8 + 16 + 64 = 88 >$. Howard has \$100 - \\$88 = \\$<100 - \\$8 = 12 > \text{left.} \#### 12
Dax went to their farm to pick some apples and found half as many bugs as ants in the garden. If there were 50 ants, calculate the total number of insects in the garden.	The total number of insects in the garden is 25 bugs + 50 ants = $^{\circ}$ 25+50=75 $^{\circ}$ 75 insects #### 75	if there were 50 ants, and there were half as many bugs as ants, there were $1/2 * 50 = 1/2 * 50 = 1 * 50 / 2 =$ iseliseliseliletilletilletilletilletille
Last night Rick killed ten wolves and 15 cougars while hunting. Today Rick killed three times as many wolves as cougars and three fewer cougars than the previous night. How many animals did Rick kill?	The total number of animals that Rick killed while hunting yesterday is 15+10 = «15+10=25»25 Today, Rick killed three fewer cougars than the previous night, a total of 15-3 = «15-3=12»12 cougars. He also killed three times as many wolves as cougars today, a total of 3*12 = «3*12=36»36 wolves. Together, his hunt today yielded 36+12 = «36+12=48»48 animals. In total, he has 48+25 = «48+25=73»73 animals from the hunt. #### 73	yesterday, he killed 10 wolves and 15 cou, so today he killed 3 times as many wolves as cou, which is isel mktimeilletiselovuovuilletov

Table 21: Example of correct and incorrect (overrefused) generations of RepBend on GSM8K.

A.9 Ablation Study Sosses

Equations 14, 15, and 16 present the loss functions used for models A1, A2 and A3 in our ablation study.

$$\mathcal{L}_{A1} = \alpha \cdot \frac{1}{N} \sum_{i=1}^{N} \max(0, \mathbf{d}_{bp}(\mathbf{h}_{b,i}, \mathbf{h}'_{b,i}) + m_b)$$

$$+ \beta \cdot \frac{1}{N} \sum_{i=1}^{N} \max(0, -\mathbf{d}_{hn}(\mathbf{h}'_{h,i}, \mathbf{h}_{h,i}) + m_h)$$

$$+ \gamma \cdot D_{KL}(\pi(b), \pi'(b))$$
(14)

$$\mathcal{L}_{A2} = \alpha \cdot \frac{1}{N} \sum_{i=1}^{N} \max(0, \mathbf{d}_{bp}(\mathbf{h}_{b,i}, \mathbf{h}'_{b,i}) + m_b)$$

$$+ \beta \cdot \frac{1}{N} \sum_{i=1}^{N} \max(0, \mathbf{d}_{hp}(\mathbf{h}'_{h,i}, \hat{\mathbf{h}}') - \mathbf{d}_{hn}(\mathbf{h}'_{h,i}, \mathbf{h}_{h,i}) + m_h)$$

$$+ \gamma \cdot D_{KL}(\pi(b), \pi'(b))$$
(15)

$$\mathcal{L}_{A3} = \alpha \cdot \frac{1}{N} \sum_{i=1}^{N} \max(0, \mathbf{d}_{bp}(\mathbf{h}_{b,i}, \mathbf{h}'_{b,i}) - \mathbf{d}_{bn}(\mathbf{h}'_{b,i}, \hat{\mathbf{h}}') + m_b)$$

$$+ \beta \cdot \frac{1}{N} \sum_{i=1}^{N} \max(0, \mathbf{d}_{hp}(\mathbf{h}'_{h,i}, \hat{\mathbf{h}}') - \mathbf{d}_{hn}(\mathbf{h}'_{h,i}, \mathbf{h}_{h,i}) + m_h)$$

$$+ \gamma \cdot D_{KL}(\pi(b), \pi'(b))$$
(16)

A.10 Proofs

Theorem 1. The circuit breakers loss \mathcal{L}_{CB} can be rewritten as a triplet loss $\mathcal{L}_{triplet}$ with null distances $d_0(x,y) = 0$.

Proof. We recall the definition of the circuit breakers loss:

$$\mathcal{L}_{CB} = \alpha \cdot \left\| \mathbf{h}_{b,i} - \mathbf{h}'_{b,i} \right\|_{2}^{2} + \beta \cdot \text{ReLU}(\cos_{-}\sin(\mathbf{h}_{h,i}, \mathbf{h}'_{h,i}))$$
(17)

$$= \alpha \cdot \text{ReLU}(\|\mathbf{h}_{b,i} - \mathbf{h}'_{b,i}\|_{2}^{2}) + \beta \cdot \text{ReLU}(\cos_{\sin}(\mathbf{h}_{h,i}, \mathbf{h}'_{h,i}))$$
(18)

$$= \alpha \cdot \text{ReLU}(\|\mathbf{h}_{b,i} - \mathbf{h}'_{b,i}\|_{2}^{2}) + \beta \cdot \text{ReLU}(-\mathsf{d}_{cos}(\mathbf{h}_{h,i}, \mathbf{h}'_{h,i}) + 1)$$
(19)

$$= \mathcal{L}_{triplet} \tag{20}$$

with the parameters $m_h = 1$, $m_b = 1$, $d_{bp} = d_2$, $d_{bn} = d_0$, $d_{hp} = d_0$, $d_{hn} = d_{cos}$.

Theorem 2. The RepBend loss \mathcal{L}_{RB} without the KL divergence term can be rewritten as a triplet loss $\mathcal{L}_{triplet}$ with a null distance $d_0(x,y) = 0$ and an arbitrary large margin m_h .

Proof. The definition of the RepBend loss, minus the KL divergence term, is:

$$\mathcal{L}_{RB} = \frac{1}{2} \cdot \left\| \mathbf{h}'_{b,i} - \mathbf{h}_{b,i} \right\|_2 \tag{21}$$

$$-\alpha \cdot \left\| \mathbf{h}_{h,i}' - \mathbf{h}_{h,i} \right\|_{2} \tag{22}$$

$$-\beta \cdot \cos_{\sin}(A) \tag{23}$$

where $\cos_{\sin}(A)$ is defined as the average cosine similarity between all pairs of vectors in A.

$$\cos_{-}\operatorname{sim}(A) = \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j=1, j \neq i}^{n} \cos_{-}\operatorname{sim}(\mathbf{a}_{i}, \mathbf{a}_{j})$$
(24)

We write μ the mean of the normalized vectors $\hat{\mathbf{a}}_i$.

$$\mu = \frac{1}{n} \sum_{i=1}^{n} \hat{\mathbf{a}}_i \tag{25}$$

Let us assume that the vectors are normalized, i.e. $||\mathbf{a}_i|| = 1$.

$$\cos_{-}\operatorname{sim}(A) = \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j=1, j \neq i}^{n} \cos_{-}\operatorname{sim}(\mathbf{a}_{i}, \mathbf{a}_{j})$$
(26)

$$= \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j=1, j \neq i}^{n} \mathbf{a}_i \cdot \mathbf{a}_j$$

$$(27)$$

$$= \frac{1}{n(n-1)} \sum_{i=1}^{n} \left(\sum_{j=1}^{n} \mathbf{a}_i \cdot \mathbf{a}_j \right) - \mathbf{a}_i \cdot \mathbf{a}_i$$
 (28)

$$= -\frac{1}{n-1} + \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j=1}^{n} \mathbf{a}_{i} \cdot \mathbf{a}_{j}$$
 (29)

$$= -\frac{1}{n-1} + \frac{1}{n(n-1)} \sum_{i=1}^{n} \mathbf{a}_{i} \cdot \left(\sum_{j=1}^{n} \mathbf{a}_{j}\right)$$
 (30)

$$= -\frac{1}{n-1} + \frac{1}{(n-1)} \sum_{i=1}^{n} \mathbf{a}_i \cdot \mu$$
 (31)

$$= -\frac{1}{n-1} - \frac{1}{(n-1)} \sum_{i=1}^{n} d_{\cos}(\mathbf{a}_i, \mu)$$
 (32)

Therefore, maximizing $\cos_{-}\sin(A)$ is equivalent to minimizing $\frac{1}{n}\sum_{i=1}^{n}d_{\cos}(\mathbf{a}_{i},\mu)$. Since d_{\cos} is the cosine distance, the case also holds when the vectors are not normalized.

Finally, we can rewrite the RepBend loss as:

$$\mathcal{L}_{RB} = \frac{1}{2} \cdot \left\| \mathbf{h}'_{b,i} - \mathbf{h}_{b,i} \right\|_2 \tag{33}$$

$$-\alpha \cdot \left\| \mathbf{h}_{h,i}' - \mathbf{h}_{h,i} \right\|_{2} \tag{34}$$

$$-\beta \cdot \frac{1}{H} \sum_{n=1}^{H} d_{\cos}(\mathbf{h}'_{h,n}, \mu)$$
(35)

$$= \frac{1}{2} \cdot \text{ReLU}(\mathbf{d}_2(\mathbf{h}'_{b,i}, \mathbf{h}_{b,i}) + d_0 - 0)$$
(36)

+ ReLU(
$$\alpha \cdot d_2(\mathbf{h}'_{h,i}, \mathbf{h}_{h,i}) - \frac{\beta}{n-1} \cdot d_{\cos}(\mathbf{h}'_{h,i}, \mu) + m_h$$
) (37)

where we assume that m_h is large enough to make the ReLU function non-zero. In practice, the RepBend loss would need to be stopped at a certain point, otherwise the loss would tend to $-\infty$. By adding a margin m_h , we can freely choose the number of iterations without worrying about divergence of the loss. Therefore, our triplet loss formulation of RepBend with margins is more practical and stable.