# Keep Safety Locked for finetuned LLM

## Anonymous ACL submission

#### Abstract

002 Fine-tuning large language models (LLMs) on additional datasets is often necessary to optimize them for specific downstream tasks. However, existing safety alignment measures, which restrict harmful behavior during inference, are insufficient to mitigate safety risks during finetuning. Alarmingly, fine-tuning with just 10 toxic sentences can make models comply with harmful instructions. We introduce Safety-Lock, a novel alignment intervention method 012 that maintains robust safety post-fine-tuning through efficient and transferable mechanisms. SafetyLock leverages our discovery that finetuned models retain similar safety-related activation representations to their base models. This insight enables us to extract what we term 017 the Meta-SafetyLock, a set of safety bias directions representing key activation patterns associated with safe responses in the original 021 model. We can then apply these directions universally to fine-tuned models to enhance their safety. By searching for activation directions across multiple token dimensions, SafetyLock achieves enhanced robustness and transferability. SafetyLock re-aligns fine-tuned models in under 0.01 seconds without additional computational cost. Our experiments demonstrate that SafetyLock can reduce the harmful instruction response rate from 60% to below 1% in toxic fine-tuned models. It surpasses traditional methods in both performance and efficiency, offering a scalable, non-invasive solution for ensuring the safety of customized LLMs. Our analysis across various fine-tuning scenarios confirms SafetyLock's robustness.

# 1 Introduction

037

040

043

Large language models (LLMs) have demonstrated increasing utility across various domains (Wei et al., 2022b,a; Weng et al., 2023; Hadar-Shoval et al., 2024), yet their potential to handle harmful queries has raised significant concerns (Carroll et al., 2023; Hendrycks et al., 2023). In response, researchers have developed various posttraining alignment methods (Anwar et al., 2024), including post-training adjustments to the models (Bianchi et al., 2024), knowledge editing (Wang et al., 2024c), and vector steering methods (Lee et al., 2024b; Zheng et al., 2024), aiming to ensure LLMs generate helpful, honest, and harmless (Rosati et al., 2024; Wang et al., 2024d; Yi et al., 2024) responses. These measures are expected to teach models to refuse harmful queries during inference (Huang et al., 2024b; Wang et al., 2024b; Raza et al., 2024; Zou et al., 2024). 044

045

046

047

051

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

078

081

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

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

086

087

090

094

101

102

103

104

105

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

132

133

134

136

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

Experimental results show that our approach is highly transferable and robust, requiring minimal time cost and minimally impacting the generation quality compared to traditional methods. First, we facilitate the efficient transfer of safety measures from base models to their fine-tuned variants, 137 including Llama-3-8B Instruct, Llama-3-70B In-138 struct, and Mistral-Large-2 123B. Second, Safety-139 Lock can be deployed without GPU resources 140 in less than 0.01 seconds (Sections 3.2 and 4.3), 141 highlighting our method's universality. Secondly, 142 SafetyLock significantly reduces the ASR from 143 54.24% to 0.03% in fine-tuned language models 144 and demonstrates robust resistance to both typical 145 safety attacks and dual attacks with prompt-based 146 methods. With the help of SafetyLock, we de-147 crease ASR from 98% to 2% for DeepInception 148 attacks (Sections 4.2 and 4.4). Finally, we con-149 ducted experiments on eight general tasks, demon-150 strating minimal performance decay. We show that 151 SafetyLock maintains a high response rate, with a 152 slight decrease from 99.4% to 98.1% (Sections 4.3 153 and 4.5). Our work advances the field of LLM 154 safety alignment by introducing Meta-SafetyLock, 155 a framework that fundamentally reimagines how 156 safety measures can be efficiently distributed across 157 fine-tuned models. While previous works estab-158 lished important foundations through safety vectors 159 (Bhardwaj et al., 2024) and various safety interven-160 tion methods (Zhao et al., 2024; Hazra et al., 2024; 161 Yi et al., 2024), our approach uniquely operates at 162 the attention-head level, supported by our discovery 163 that safety-relevant attention heads maintain consis-164 tency even after fine-tuning. This insight enables us 165 to extract a single Meta-SafetyLock from the base 166 model that can be rapidly deployed across multi-167 ple fine-tuned variants without requiring repeated 168 safety pattern searches, achieving remarkable effi-169 ciency without GPU resources. 170

# 2 Related Work

Alignment of LLMs. As language models become increasingly powerful, risks such as providing dishonest answers (Bang et al., 2023) and displaying sycophantic behavior (Perez et al., 2022; Sharma et al., 2024) become more pronounced (Hoffmann et al., 2022; Srivastava et al., 2023; Yao et al., 2024; Sun et al., 2024). Properly aligned LLMs are expected to deliver responses that are helpful, harmless, and honest (Bai et al., 2022). Specifically, harmlessness is addressed through safety alignment (Ji et al., 2024; Zhao et al., 2024), which involves equipping LLMs with safety protocols that enable them to decline harmful instructions. Common approaches for safety alignment include instruction tuning (Ouyang et al., 2022; Zhang et al., 2024b), 171

172

173

174

175

176

177

178

179

180

181

182

183

184

185



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

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

187

188

190

191 192

193

194

197

198

199

201

210

211

212

214

Safeguards of LLMs. Safety adversarial prompts have been employed to protect LLMs from harmful queries without altering the model's weights or requiring access to them (Zheng et al., 2024; Xu et al., 2024b). These prompts are added to the system prompt text to defend against jailbreak attacks (Shi et al., 2023; Hong et al., 2024). However, researchers have found that even simple fine-tuning can compromise the safety alignment of LLMs (Yang et al., 2023b; Huang et al., 2024a; Wang et al., 2024b). For example, Qi et al. (2024) demonstrated that using just 10 harmful examples was sufficient to undermine the safety alignment of GPT-3.5. This finding underscores the lack of robustness in current safety alignment strategies. Recent works have made progress in understanding safety mechanisms - from identifying safety neurons (Chen et al., 2024b) to revealing the role of feed-forward layers in safety responses (Geva et al.,

2021). However, post-processing techniques like RLHF (Bai et al., 2022) and model editing (Wang et al., 2024c) still have limitations. For instance, PPO and DPO adjust the entire activation space, while model editing targets concentrated areas, often missing dispersed safety information.

215

216

217

218

219

220

221

223

224

225

226

227

228

229

230

231

232

233

235

236

237

238

Interventions in LLMs. Intervening in the internal activation of Transformer-based language models during inference can trigger specific transformations (Wu et al., 2024b; Turner et al., 2023; Rimsky et al., 2023). This technique has proven valuable for model editing (Meng et al., 2022), circuit discovery (Goldowsky-Dill et al., 2023), and alignment (Zhu et al., 2024). Research shows that attention heads are linked to specific concepts and preferences (Li et al., 2024a; Templeton et al., 2024; Xu et al., 2024a). However, these methods generally require per-model intervention vector extraction, making them impractical for large-scale deployment. Building on this, SafetyLock achieves precise safety alignment through multi-token-level interventions, using only the activation values from the original model, thus providing robustness to parameter changes while enhancing efficiency.

## **3** Method: SafetyLock

As illustrated in Figure 1b, SafetyLock comprises 240 two main phases: manufacturing Meta-SafetyLock 241 and distributing SafetyLock. Our research reveals that fine-tuned LLMs and their original counterparts exhibit remarkably similar activation representations and intervention directions for safety preferences (Section A). This discovery enables us to extract safety bias directions (Meta-SafetyLock) from the original LLMs (Section 3.2). To enhance universality and robustness, we propose searching for activation directions across multiple tokens, thereby improving the transferability of SafetyLock during distribution (Section 3.3).

## 3.1 Safety features in Transformer

242

243

244

246

247

251

254

255

262

265

267

268

270

272

273

274

275

278

281

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

## 3.2 Constructing Meta-SafetyLock

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

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

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

291

292

293

294

295

299

301

302

303

304

306

307

308

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

329

331

332

333

334

335

337

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

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

### 3.3 Distributing SafetyLock

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

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

shift being proportional to the standard deviation, 338 scaled by a factor  $\alpha$ .

**Experiments** 

4

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

$$\operatorname{Bias}_{l} = \operatorname{Bias}_{l} + \alpha \sum_{h=1}^{H} Q_{l}^{h} \left( \sigma_{l}^{h} \theta_{l}^{h} \right). \quad (3)$$

In this section, we present experiments to evaluate

the effectiveness of the SafetyLock in enhancing

model safety and inference efficiency, while main-

Threat Model Selections. Following previous red

teaming and safeguarding studies on aligned LLMs

(Yuan et al., 2024), we consider a threat model

where attackers can fine-tune aligned LLMs, typi-

cally through API access to closed-source models.

The primary objective is jailbreaking these models

and removing safety constraints (Wei et al., 2023;

Carlini et al., 2023) while SafetyLock aims to re-

build the safety guard. We use Llama-3-8B Chat,

Llama-3-70B Chat, and Mistral-Large-2 123B as

our base models, fine-tuning them on datasets rep-

resenting each risk level to simulate real-world sce-

narios. Please refer to Appendix D for detailed

ments on three risks: (1) explicitly harmful datasets,

where attackers intentionally fine-tune models on

malicious content (Ganguli et al., 2022; Qi et al.,

2023); (2) implicitly harmful datasets, which may

appear benign but lead to compromised safety

guardrails (Qi et al., 2024); and (3) benign datasets,

where even well-intentioned fine-tuning can inad-

Fine-tuning Datasets. We conducted experi-

taining model's general performance.

4.1 Experimental Details

baseline experimental setups.

343

342

341

- 345

347

- 349

361 362

364

371

373 374

378

381

384

vertently degrade model safety (Wang et al., 2023b). For Risk-1, we use negative samples from the HH-RLHF dataset (Bai et al., 2022). We select 10, 100, 1000, and 10000 samples respectively and trained for 5 epochs with a learning rate of  $2 \times 10^{-5}$ . For Risk-2, we use 10 samples from Qi et al. (2024) and train for 5 epochs with a learning rate of  $2 \times 10^{-5}$ . For Risk-3, we used the first 50,000 samples from the Alpaca dataset (Wang et al., 2023b) and trained for 5 epochs with a learning rate of  $2 \times 10^{-5}$ .

Safety Evaluation and Metrics. Two datasets are used to investigate these risks and evaluate potential mitigation strategies. HEx-PHI (Qi et al.,

2024) is based on 11 categories of prohibited use cases merged from Meta's Llama-3 acceptable use policy and OpenAI's usage policies. The dataset includes 30 examples per category, totalling 330 examples. This ensures a comprehensive safety evaluation aligned with industry-standard usage policies. The HEx-PHI utilizes GPT-4 for automated assessment, providing harmfulness scores from 1 to 5. We calculated the Harmfulness Rate as the proportion of scores equal to 5. AdvBench is released by (Zou et al., 2023b), we adhere to the original paper's setup and calculate the ASR through string matching.

385

386

387

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

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

## 4.2 Results over Different Risk Levels

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

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

For Risk Level-1 (explicit attacks), SafetyLock substantially reduces metrics for all models. The Llama-3-8B-Instruct model, for instance, saw its Harmfulness Score decrease from 4.13 to 1.36, Rate from 70.01% to 3.33%, and ASR from 49.24% to 0.19%. Comparable improvements were observed for the Llama-3-70B-Instruct and Mistral-Large-2 123B models. Risk Level-2 and Risk Level-3 also showed significant improve-

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

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



Figure 2: Impact of increasing harmful training samples on model safety with and without SafetyLock.

ments. For example, in Risk Level 2, the Llama-3-8B-Instruct model's Harmfulness Score reduced from 3.19 to 1.07. Similar improvements were observed across all model sizes, demonstrating Safety-Lock's ability to maintain ethical guardrails during routine model customization processes.

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

In Figure 2, we further supplement an ablation with larger training sets on risk 1 (100, 1000, and 10000 harmful samples) showing that SafetyLockprotected models maintain low ASR across all sample sizes. Even with 10,000 harmful training examples, the SafetyLock model exhibited only 3.46% ASR, compared to 62.31% for the unprotected model. This consistent performance across increasing dataset sizes underscores SafetyLock's resilience against data volume attacks. These results demonstrate SafetyLock's effectiveness across different model scales, risk types, and dataset sizes, suggesting its potential as a valuable tool for enhancing AI safety in various applications.

## 4.3 Comparative Analysis

To comprehensively evaluate SafetyLock's efficacy,
we conducted a comparative analysis against established baseline methods, categorized into trainingbased and inference-time approaches, as illustrated
in Figure 3. This analytical framework enables a
thorough assessment of various strategies for main-

taining model safety in fine-tuned language models.

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

As demonstrated in Figure 3, in terms of efficiency, SafetyLock exhibits a remarkable computational economy. Its inference time of 0.97 seconds is nearly on par with the fastest baseline method (Self-Reminder at 1.12 seconds), while its training time of 0.01 seconds and additional GPU memory usage of 0.0 GB are orders of magnitude lower than all training-based methods. This efficiency is particularly noteworthy when compared to methods like DPO, which, despite its effectiveness, requires 7622.0 seconds of training time and 45.12 GB of GPU memory. Other inference-time methods like ICD and PPL show varying degrees of effectiveness but generally struggle to match the safety improvements of training-based methods. SFT with safety data mixing post-fine-tuning offers a more balanced approach, achieving a Harmfulness Score of 1.03 with reduce resource requirements of 779 seconds and 38.32 GB GPU memory. Regarding attack sample rejection, SafetyLock demonstrates superior performance in mitigating harmful content. It achieves a Harmfulness Score of 1.04, equivalent to that achieved by models undergoing safety realignment via DPO, indicating its exceptional ability to reduce the generation of harmful content. Furthermore, SafetyLock's AdvBench ASR of 0.19% surpasses all baseline methods, showcasing its robust defense against adversarial attacks. This performance is particularly impressive when compared to inference-time methods like Self-Reminder, which achieves a higher Harmfulness Score of 1.82 and an AdvBench ASR of 19.81%.

We further assess the models' performance on benign inputs to ensure safety enhancements did not compromise normal text generation by selecting 500 test samples from the Alpaca dataset. The results reveal that SafetyLock preserves a 98.1% normal response rate, closely trailing the origi-



Figure 3: Comparison of Methods for Mitigating Safety Risks in Fine-tuned Language Models (Llama-3-Instruct 8B). Upper row: Compared with inference-time methods; Lower row: Compared with training-time methods.

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

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

nal Vanilla model's 99.4%. Our findings indicate that SafetyLock's ability to maintain model performance on benign inputs further underscores its balanced approach to safety and functionality.

502

505

506

508

509

510

511

513

514

515

516

517

518

519

521

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

# 4.4 Against Combined Attacks

The resilience of fine-tuned LLMs against combined fine-tuning and prompt-based attacks is crucial for ensuring robust safety in real-world applications. To further assess robustness, we introduced a combined attack scenario: fine-tuning model attacks followed by prompt-based attacks. We evaluated four commonly prompt attack methods: AutoDAN (Liu et al., 2024), DeepInception (Li et al., 2024b), GCG (Zou et al., 2023b), PAIR (Chao et al., 2024), and XSTest (Röttger et al., 2023) comparing their performance against several defense techniques, as illustrated in Table 2. 522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

SafetyLock demonstrates exceptional effectiveness across all tested attack methods. For Auto-DAN attacks, SafetyLock reduces the ASR to a mere 4.0%, significantly outperforming other methods such as ICD (46.0%) and Self-Exam (66.0%). Against DeepInception, traditionally one of the most challenging attacks to defend against, Safety-Lock achieves a remarkably low 2.0% ASR, while all other methods fail to provide any meaningful defense (98.0% ASR across the board). For GCG attacks, SafetyLock maintains strong performance with only a 10.0% ASR, second only to PPL's 0.0% but considerably better than most other methods, including Vanilla (74.0%) and Retokenization (94.0%). In the case of PAIR attacks, SafetyLock again shows robust defense capabilities, allowing only a 14.0% ASR, outperforming all other tested methods. Additionally, on the structured XSTest



Figure 4: Performance comparison of various methods on downstream tasks.

benchmark, SafetyLock achieves a state-of-the-art 4.0% ASR, substantially outperforming other approaches such as ICD (7.0%) and Self-Reminder (8.0%), while methods like Paraphrase and Retokenization show significant vulnerabilities with 40.0% and 57.5% ASR respectively.

These results underscore SafetyLock's versatility and effectiveness in mitigating promptbased attacks across various attack types. Its consistent performance demonstrates a comprehensive approach to model safety, addressing the complex challenges posed by diverse attack scenarios in language model deployment. The ability to maintain such low ASR across different attack methods suggests that SafetyLock provides a more generalizable and robust defense mechanism.

### 4.5 Generalization Capabilities of SafetyLock

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

As illustrated in Figure 4, SafetyLock demonstrates remarkable ability to maintain model performance while ensuring safety. Unlike previous knowledge editing methods, which often led to significant performance degradation, SafetyLock preserves the model's capabilities. For instance, on the AddSub task, SafetyLock maintains 85.57% performance (compared to original 86.33%), while Model-Edited shows complete performance collapse. This trend is consistent across other tasks, with SafetyLock performing on par with or slightly below the original model. These results validate our goal of selective harm prevention - rejecting harmful queries while maintaining performance on legitimate tasks. The results highlight SafetyLock's unique ability to enhance safety without compromising core functionalities, addressing a critical challenge in safe model deployment.

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

## 5 Conclusion

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

583

584

587

553

# 621 Limitations

While SafetyLock demonstrates promising results in maintaining the safety of fine-tuned language 623 models, it is important to acknowledge several lim-624 itations. Primarily, SafetyLock requires access to 625 both model weights and intermediate activations for implementation, which may limit its applica-627 bility in scenarios where such access is restricted or unavailable. Additionally, the method employs a symmetric locking mechanism; consequently, if an unauthorized party gains access to the model weights or activation values, they could potentially reverse-engineer the process to unlock and bypass SafetyLock's protections. Lastly, while SafetyLock shows strong performance against current attack methods, its long-term robustness against evolving 636 adversarial techniques remains to be studied. These 637 limitations present opportunities for future work to enhance and expand the capabilities of SafetyLock, ensuring its continued effectiveness in maintaining AI safety. 641

#### References

642

647

649

654

663

668

671

672

- Gabriel Alon and Michael Kamfonas. 2023. Detecting language model attacks with perplexity. *Preprint*, arXiv:2308.14132.
- Usman Anwar, Abulhair Saparov, Javier Rando, Daniel Paleka, Miles Turpin, Peter Hase, Ekdeep Singh Lubana, Erik Jenner, Stephen Casper, Oliver Sourbut, et al. 2024. Foundational challenges in assuring alignment and safety of large language models. *arXiv preprint arXiv:2404.09932*.
- Patel Arkil, Bhattamishra Satwik, and Goyal Navin. 2021. Are nlp models really able to solve simple math word problems?
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv*:2204.05862.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. In *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 675–718.
- Rishabh Bhardwaj, Do Duc Anh, and Soujanya Poria. 2024. Language models are homer simpson! safety

re-alignment of fine-tuned language models through task arithmetic. *arXiv preprint arXiv:2402.11746*.

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

- Manish Bhatt, Sahana Chennabasappa, Cyrus Nikolaidis, Shengye Wan, Ivan Evtimov, Dominik Gabi, Daniel Song, Faizan Ahmad, Cornelius Aschermann, Lorenzo Fontana, Sasha Frolov, Ravi Prakash Giri, Dhaval Kapil, Yiannis Kozyrakis, David LeBlanc, James Milazzo, Aleksandar Straumann, Gabriel Synnaeve, Varun Vontimitta, Spencer Whitman, and Joshua Saxe. 2023. Purple Ilama cyberseceval: A secure coding benchmark for language models. *Preprint*, arXiv:2312.04724.
- Federico Bianchi, Mirac Suzgun, Giuseppe Attanasio, Paul Rottger, Dan Jurafsky, Tatsunori Hashimoto, and James Zou. 2024. Safety-tuned LLaMAs: Lessons from improving the safety of large language models that follow instructions. In *The Twelfth International Conference on Learning Representations*.
- Nicholas Carlini, Milad Nasr, Christopher A Choquette-Choo, Matthew Jagielski, Irena Gao, Anas Awadalla, Pang Wei Koh, Daphne Ippolito, Katherine Lee, Florian Tramer, et al. 2023. Are aligned neural networks adversarially aligned? *arXiv preprint arXiv:2306.15447*.
- Micah Carroll, Alan Chan, Henry Ashton, and David Krueger. 2023. Characterizing manipulation from ai systems. *Preprint*, arXiv:2303.09387.
- 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.
- Canyu Chen, Baixiang Huang, Zekun Li, Zhaorun Chen, Shiyang Lai, Xiongxiao Xu, Jia-Chen Gu, Jindong Gu, Huaxiu Yao, Chaowei Xiao, Xifeng Yan, William Yang Wang, Philip Torr, Dawn Song, and Kai Shu. 2024a. Can editing llms inject harm? *arXiv preprint arXiv: 2407.20224*.
- Jianhui Chen, Xiaozhi Wang, Zijun Yao, Yushi Bai, Lei Hou, and Juanzi Li. 2024b. Finding safety neurons in large language models. *Preprint*, arXiv:2406.14144.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpacafarm: A simulation framework for methods that learn from human feedback. *Preprint*, arXiv:2305.14387.
- Pranav M. Gade, Simon Lermen, Charlie Rogers-Smith, and Jeffrey Ladish. 2023. Badllama: cheaply removing safety fine-tuning from llama 2-chat 13b. *ArXiv*, abs/2311.00117.

839

Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda 727 Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, 728 Ethan Perez, Nicholas Schiefer, Kamal Ndousse, Andy Jones, Sam Bowman, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Nelson Elhage, Sheer El-Showk, Stanislav Fort, Zac Hatfield-Dodds, Tom Henighan, Danny Hernandez, Tristan Hume, Josh Jacobson, Scott Johnston, Shauna Kravec, Catherine Olsson, Sam Ringer, Eli Tran-Johnson, 735 Dario Amodei, Tom Brown, Nicholas Joseph, Sam McCandlish, Chris Olah, Jared Kaplan, and Jack Clark. 2022. Red teaming language models to re-738 duce harms: Methods, scaling behaviors, and lessons learned. Preprint, arXiv:2209.07858.

731

737

741

742

743

744 745

746

747

748

750

751

753

754

755

756

757

761

763

764

766

767

768

770

774

775

776

777

778

779

- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. Transformer feed-forward layers are key-value memories. Preprint, arXiv:2012.14913.
  - Charles Godfrey, Davis Brown, Tegan Emerson, and Henry Kvinge. 2022. On the symmetries of deep learning models and their internal representations. In Advances in Neural Information Processing Systems, volume 35, pages 11893–11905. Curran Associates, Inc.
  - Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. 2023. Localizing model behavior with path patching. arXiv preprint arXiv:2304.05969.
  - Satvik Golechha and James Dao. 2024. Challenges in mechanistically interpreting model representations. Preprint, arXiv:2402.03855.
  - Dorit Hadar-Shoval, Kfir Asraf, Yonathan Mizrachi, Yuval Haber, and Zohar Elyoseph. 2024. Assessing the alignment of large language models with human values for mental health integration: Cross-sectional study using schwartz's theory of basic values. JMIR Mental Health, 11:e55988.
  - Rima Hazra, Sayan Layek, Somnath Banerjee, and Soujanya Poria. 2024. Safety arithmetic: A framework for test-time safety alignment of language models by steering parameters and activations. arXiv preprint arXiv:2406.11801.
- Dan Hendrycks, Geoffrey Hinton, Yoshua Bengio, Demis Hassabis, Sam Altman, Dario Amodei, Dawn Song, Ted Lieu, Bill Gates, Ya-Qin Zhang, Ilya Sutskever, Igor Babuschkin, Shane Legg, Martin Hellman, James Manyika, Yi Zeng, and Xianyuan Zhan. 2023. Statement on ai risk. https:// www.safe.ai/work/statement-on-ai-risk. Accessed: 2024-06-20.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Thomas Hennigan, Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karén Simonyan, Erich Elsen, Oriol Vinyals, Jack Rae, and Laurent Sifre. 2022. An empirical analysis

of compute-optimal large language model training. In Advances in Neural Information Processing Systems, volume 35, pages 30016–30030. Curran Associates, Inc.

- Zhang-Wei Hong, Idan Shenfeld, Tsun-Hsuan Johnson Wang, Yung-Sung Chuang, Aldo Pareja, James Glass, Akash Srivastava, and Pulkit Agrawal. 2024. Curiosity-driven red-teaming for large language models. ArXiv, abs/2402.19464.
- Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. 2014. Learning to solve arithmetic word problems with verb categorization. empirical methods in natural language processing.
- Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. 2024a. Lazy safety alignment for large language models against harmful finetuning.
- Youcheng Huang, Jingkun Tang, Duanyu Feng, Zheng Zhang, Wenqiang Lei, Jiancheng Lv, and Anthony G Cohn. 2024b. Dishonesty in helpful and harmless alignment. arXiv preprint arXiv:2406.01931.
- 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.
- Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2024. Beavertails: Towards improved safety alignment of llm via a humanpreference dataset. Advances in Neural Information Processing Systems, 36.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. In Advances in Neural Information Processing Systems.
- Andrew Lee, Xiaoyan Bai, Itamar Pres, Martin Wattenberg, Jonathan K Kummerfeld, and Rada Mihalcea. 2024a. A mechanistic understanding of alignment algorithms: A case study on dpo and toxicity. arXiv preprint arXiv:2401.01967.
- Bruce W. Lee, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Erik Miehling, Pierre Dognin, Manish Nagireddy, and Amit Dhurandhar. 2024b. Programming refusal with conditional activation steering. *Preprint*, arXiv:2409.05907.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2024a. Inferencetime intervention: Eliciting truthful answers from a language model. Advances in Neural Information Processing Systems, 36.
- Xuan Li, Zhanke Zhou, Jianing Zhu, Jiangchao Yao, Tongliang Liu, and Bo Han. 2024b. Deepinception: Hypnotize large language model to be jailbreaker. Preprint, arXiv:2311.03191.

Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. 2024. AutoDAN: Generating stealthy jailbreak prompts on aligned large language models. In *The Twelfth International Conference on Learning Representations*.

841

849

853

854

855

857

859

864

867

868

870

871

872

873

874

875

876

877

883

884

892 893

896

- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372.
- Yu Meng, Mengzhou Xia, and Danqi Chen. 2024. Simpo: Simple preference optimization with a reference-free reward. *arXiv preprint arXiv:2405.14734*.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. 2021. Fast model editing at scale. *arXiv preprint arXiv:2110.11309*.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. 2022. Memorybased model editing at scale. In *International Conference on Machine Learning*, pages 15817–15831. PMLR.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- ShengYun Peng, Pin-Yu Chen, Matthew Hull, and Duen Horng Chau. 2024. Navigating the safety landscape: Measuring risks in finetuning large language models. *Preprint*, arXiv:2405.17374.
- Ethan Perez, Sam Ringer, Kamilė Lukošiūtė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Ben Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jackson Kernion, James Landis, Jamie Kerr, Jared Mueller, Jeeyoon Hyun, Joshua Landau, Kamal Ndousse, Landon Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav Kingsland, Nelson Elhage, Nicholas Joseph, Noemí Mercado, Nova DasSarma, Oliver Rausch, Robin Larson, Sam McCandlish, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Jack Clark, Samuel R. Bowman, Amanda Askell, Roger Grosse, Danny Hernandez, Deep Ganguli, Evan Hubinger, Nicholas Schiefer, and Jared Kaplan. 2022. Discovering language model behaviors with model-written evaluations. Preprint, arXiv:2212.09251.
- Mansi Phute, Alec Helbling, Matthew Daniel Hull, ShengYun Peng, Sebastian Szyller, Cory Cornelius, and Duen Horng Chau. 2024. LLM self defense: By

self examination, LLMs know they are being tricked. In *The Second Tiny Papers Track at ICLR 2024*.

- Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Peter Henderson, Mengdi Wang, and Prateek Mittal. 2023. Visual adversarial examples jailbreak aligned large language models. *Preprint*, arXiv:2306.13213.
- Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. 2024. Finetuning aligned language models compromises safety, even when users do not intend to! In *The Twelfth International Conference on Learning Representations*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.
- Shaina Raza, Oluwanifemi Bamgbose, Shardul Ghuge, Fatemeh Tavakoli, and Deepak John Reji. 2024. Developing safe and responsible large language models–a comprehensive framework. *arXiv preprint arXiv:2404.01399*.
- Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner. 2023. Steering Ilama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*.
- Domenic Rosati, Jan Wehner, Kai Williams, Lukasz Bartoszcze, Jan Batzner, Hassan Sajjad, and Frank Rudzicz. 2024. Immunization against harmful finetuning attacks. In *Conference on Empirical Methods in Natural Language Processing*.
- Paul Röttger, Hannah Rose Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy. 2023. Xstest: A test suite for identifying exaggerated safety behaviours in large language models. *arXiv preprint arXiv:2308.01263*.
- Subhro Roy and Dan Roth. 2016. Solving general arithmetic word problems. *arXiv: Computation and Language*.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *Preprint*, arXiv:1707.06347.
- Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R. Bowman, Esin DURMUS, Zac Hatfield-Dodds, Scott R Johnston, Shauna M Kravec, Timothy Maxwell, Sam Mc-Candlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang, and Ethan Perez. 2024. Towards understanding sycophancy in language models. In *The Twelfth International Conference on Learning Representations*.
- Zhouxing Shi, Yihan Wang, Fan Yin, Xiangning Chen, Kai-Wei Chang, and Cho-Jui Hsieh. 2023. Red teaming language model detectors with language models. *Transactions of the Association for Computational Linguistics*, 12:174–189.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, 953 Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, 957 Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea 962 Madotto, Andrea Santilli, Andreas Stuhlmüller, An-963 964 drew Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, 965 Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia 970 Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, 971 Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk 973 Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan 974 Orinion, Cameron Diao, Cameron Dour, Cather-975 ine Stinson, Cedrick Argueta, César Ferri Ramírez, 976 Chandan Singh, Charles Rathkopf, Chenlin Meng, 977 978 Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí 985 González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Do-987 han, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina 991 Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, El-993 lie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, 995 Eric Tang, Erkut Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice En-997 gefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, 998 Fatemeh Siar, Fernando Martínez-Plumed, Francesca 999 Happé, Francois Chollet, Frieda Rong, Gaurav 1000 Mishra, Genta Indra Winata, Gerard de Melo, Ger-1001 mán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovitch-1002 López, Gregor Betz, Guy Gur-Ari, Hana Galijase-1003 1004 vic, Hannah Kim, Hannah Rashkin, Hannaneh Ha-1005 jishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, 1006 Hinrich Schütze, Hiromu Yakura, Hongming Zhang, 1007 Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, 1008 Jack Geissinger, Jackson Kernion, Jacob Hilton, Jae-1009 hoon Lee, Jaime Fernández Fisac, James B. Simon, 1010 James Koppel, James Zheng, James Zou, Jan Kocoń, Jana Thompson, Janelle Wingfield, Jared Kaplan, 1011 Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, 1012 Jason Wei, Jason Yosinski, Jekaterina Novikova, 1013 Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen 1014 1015 Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Ji-1016 aming Song, Jillian Tang, Joan Waweru, John Bur-

den, John Miller, John U. Balis, Jonathan Batchelder, 1017 Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose 1018 Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, 1019 Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, 1020 Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl 1021 Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh D. Dhole, Kevin Gim-1023 pel, Kevin Omondi, Kory Mathewson, Kristen Chi-1024 afullo, Ksenia Shkaruta, Kumar Shridhar, Kyle Mc-1025 Donell, Kyle Richardson, Laria Reynolds, Leo Gao, 1026 Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-1027 Ochando, Louis-Philippe Morency, Luca Moschella, 1028 Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng 1029 He, Luis Oliveros Colón, Luke Metz, Lütfi Kerem 1030 Şenel, Maarten Bosma, Maarten Sap, Maartje ter 1031 Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas 1032 Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramírez Quintana, Marie Tolkiehn, 1034 Mario Giulianelli, Martha Lewis, Martin Potthast, 1035 Matthew L. Leavitt, Matthias Hagen, Mátyás Schu-1036 bert, Medina Orduna Baitemirova, Melody Arnaud, 1037 Melvin McElrath, Michael A. Yee, Michael Co-1038 hen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michał Swędrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike 1041 Cain, Mimee Xu, Mirac Suzgun, Mitch Walker, 1042 Mo Tiwari, Mohit Bansal, Moin Aminnaseri, Mor Geva, Mozhdeh Gheini, Mukund Varma T, Nanyun 1044 Peng, Nathan A. Chi, Nayeon Lee, Neta Gur-Ari 1045 Krakover, Nicholas Cameron, Nicholas Roberts, 1046 Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas 1047 Deckers, Niklas Muennighoff, Nitish Shirish Keskar, 1048 Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan 1049 Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, 1050 Omer Levy, Owain Evans, Pablo Antonio Moreno 1051 Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, 1052 Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, 1053 Percy Liang, Peter Chang, Peter Eckersley, Phu Mon 1054 Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, 1055 Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing 1056 Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta 1057 Rudolph, Raefer Gabriel, Rahel Habacker, Ramon 1058 Risco, Raphaël Millière, Rhythm Garg, Richard 1059 Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman 1061 Novak, Roman Sitelew, Ronan LeBras, Rosanne 1062 Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhut-1063 dinov, Ryan Chi, Ryan Lee, Ryan Stovall, Ryan 1064 Teehan, Rylan Yang, Sahib Singh, Saif M. Moham-1065 mad, Sajant Anand, Sam Dillavou, Sam Shleifer, 1066 Sam Wiseman, Samuel Gruetter, Samuel R. Bow-1067 man, Samuel S. Schoenholz, Sanghyun Han, San-1068 jeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan 1069 Ghosh, Sean Casey, Sebastian Bischoff, Sebastian 1070 Gehrmann, Sebastian Schuster, Sepideh Sadeghi, 1071 Shadi Hamdan, Sharon Zhou, Shashank Srivastava, 1072 Sherry Shi, Shikhar Singh, Shima Asaadi, Shixi-1073 ang Shane Gu, Shubh Pachchigar, Shubham Tosh-1074 niwal, Shyam Upadhyay, Shyamolima, Debnath, 1075 Siamak Shakeri, Simon Thormeyer, Simone Melzi, 1076 Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, 1077 Spencer Torene, Sriharsha Hatwar, Stanislas De-1078 haene, Stefan Divic, Stefano Ermon, Stella Bider-1079 man, Stephanie Lin, Stephen Prasad, Steven T. Pi-1080

antadosi, Stuart M. Shieber, Summer Misherghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. Preprint, arXiv:2206.04615.

1081

1082

1083

1085

1089

1090

1091

1092

1093

1096

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114 1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. In Advances in Neural Information Processing Systems, volume 33, pages 3008–3021. Curran Associates, Inc.
- Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Yuan Li, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, Zhengliang Liu, Yixin Liu, Yijue Wang, Zhikun Zhang, Bertie Vidgen, Bhavya Kailkhura, Caiming Xiong, Chaowei Xiao, Chunyuan Li, Eric Xing, Furong Huang, Hao Liu, Heng Ji, Hongyi Wang, Huan Zhang, Huaxiu Yao, Manolis Kellis, Marinka Zitnik, Meng Jiang, Mohit Bansal, James Zou, Jian Pei, Jian Liu, Jianfeng Gao, Jiawei Han, Jieyu Zhao, Jiliang Tang, Jindong Wang, Joaquin Vanschoren, John Mitchell, Kai Shu, Kaidi Xu, Kai-Wei Chang, Lifang He, Lifu Huang, Michael Backes, Neil Zhenqiang Gong, Philip S. Yu, Pin-Yu Chen, Quanguan Gu, Ran Xu, Rex Ying, Shuiwang Ji, Suman Jana, Tianlong Chen, Tianming Liu, Tianyi Zhou, William Wang, Xiang Li, Xiangliang Zhang, Xiao Wang, Xing Xie, Xun Chen, Xuyu Wang, Yan Liu, Yanfang Ye, Yinzhi Cao, Yong Chen, and Yue Zhao. 2024. Trustllm: Trustworthiness in large language models. Preprint, arXiv:2401.05561.
  - Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. 2022. Challenging big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*.
    - Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2018. Commonsenseqa: A question answering challenge targeting commonsense knowledge. *north american chapter of the association for computational linguistics*.

Adly Templeton, Tom Conerly, Jonathan Marcus, Jack 1140 Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, 1141 Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy 1142 Cunningham, Nicholas L Turner, Callum McDougall, 1143 Monte MacDiarmid, Alex Tamkin, Esin Durmus, 1144 Tristan Hume, Francesco Mosconi, C. Daniel Free-1145 man, Theodore R. Sumers, Edward Rees, Joshua 1146 Batson, Adam Jermyn, Shan Carter, Chris Olah, and 1147 Tom Henighan. 2024. Scaling monosemanticity: Ex-1148 tracting interpretable features from claude 3 sonnet. 1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

- Alex Turner, Lisa Thiergart, David Udell, Gavin Leech, Ulisse Mini, and Monte MacDiarmid. 2023. Activation addition: Steering language models without optimization. *arXiv preprint arXiv:2308.10248*.
- Guanchu Wang, Yu-Neng Chuang, Ruixiang Tang, Shaochen Zhong, Jiayi Yuan, Hongye Jin, Zirui Liu, Vipin Chaudhary, Shuai Xu, James Caverlee, and Xia Hu. 2024a. Taylor unswift: Secured weight release for large language models via Taylor expansion. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6928–6941, Miami, Florida, USA. Association for Computational Linguistics.
- Jiong Wang, Jiazhao Li, Yiquan Li, Xiangyu Qi, Junjie Hu, Yixuan Li, Patrick Drew McDaniel, Muhao Chen, Bo Li, and Chaowei Xiao. 2024b. Mitigating fine-tuning jailbreak attack with backdoor enhanced alignment. *ArXiv*, abs/2402.14968.
- Mengru Wang, Ningyu Zhang, Ziwen Xu, Zekun Xi, Shumin Deng, Yunzhi Yao, Qishen Zhang, Linyi Yang, Jindong Wang, and Huajun Chen. 2024c. Detoxifying large language models via knowledge editing. *Preprint*, arXiv:2403.14472.
- Peng Wang, Ningyu Zhang, Xin Xie, Yunzhi Yao, Bozhong Tian, Mengru Wang, Zekun Xi, Siyuan Cheng, Kangwei Liu, Guozhou Zheng, et al. 2023a. Easyedit: An easy-to-use knowledge editing framework for large language models. *arXiv preprint arXiv:2308.07269*.
- Pengyu Wang, Dong Zhang, Linyang Li, Chenkun Tan, Xinghao Wang, Ke Ren, Botian Jiang, and Xipeng Qiu. 2024d. Inferaligner: Inference-time alignment for harmlessness through cross-model guidance. *arXiv preprint arXiv:2401.11206*.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023b. Self-instruct: Aligning language models with self-generated instructions. *Preprint*, arXiv:2212.10560.
- Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2023. Jailbroken: How does llm safety training fail? *arXiv preprint arXiv:2307.02483*.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022a. Emergent abilities of large language models. *Transactions on Machine Learning Research*.

- 1197 1198 1199 1201 1203 1204 1205 1206
- 1207 1208 1209 1210 1211 1212 1213 1214 1215 1216 1217 1218 1219 1220 1221 1222
- 1223 1224 1225 1226 1227 1228 1229
- 1230 1231 1232 1233
- 1234 1235 1236
- 1237 1238
- 1240

- 1241 1242
- 1243 1244

1245

1246 1247

1248

1249 1250

1251 1252

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022b. Chain of thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems.
- Zeming Wei, Yifei Wang, Ang Li, Yichuan Mo, and Yisen Wang. 2024. Jailbreak and guard aligned language models with only few in-context demonstrations. Preprint, arXiv:2310.06387.
- Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Kang Liu, and Jun Zhao. 2024. Mastering symbolic operations: Augmenting language models with compiled neural networks. In The Twelfth International *Conference on Learning Representations.*
- Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao. 2023. Large language models are better reasoners with self-verification. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 2550-2575, Singapore. Association for Computational Linguistics.
- Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D Manning, and Christopher Potts. 2024a. Reft: Representation finetuning for language models. arXiv preprint arXiv:2404.03592.
- Zhengxuan Wu, Atticus Geiger, Aryaman Arora, Jing Huang, Zheng Wang, Noah D. Goodman, Christopher D. Manning, and Christopher Potts. 2024b. pyvene: A library for understanding and improving PyTorch models via interventions.
- Yueqi Xie, Jingwei Yi, Jiawei Shao, Justin Curl, Lingjuan Lyu, Qifeng Chen, Xing Xie, and Fangzhao Wu. 2023. Defending chatgpt against jailbreak attack via self-reminders. Nature Machine Intelligence, 5:1486-1496.
- Zhihao Xu, Ruixuan Huang, Xiting Wang, Fangzhao Wu, Jing Yao, and Xing Xie. 2024a. Uncovering safety risks in open-source llms through concept activation vector. arXiv preprint arXiv:2404.12038.
- Zihao Xu, Yi Liu, Gelei Deng, Yuekang Li, and Stjepan Picek. 2024b. Llm jailbreak attack versus defense techniques-a comprehensive study. arXiv preprint arXiv:2402.13457.
- J. Yang et al. 2023a. Red teaming language models via activation engineering. LessWrong.
- Xianjun Yang, Xiao Wang, Qi Zhang, Linda Ruth Petzold, William Yang Wang, Xun Zhao, and Dahua Lin. 2023b. Shadow alignment: The ease of subverting safely-aligned language models. ArXiv, abs/2310.02949.
- Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. 2024. A survey on large language model (llm) security and privacy: The good, the bad, and the ugly. High-Confidence Computing, 4(2):100211.

Xin Yi, Shunfan Zheng, Linlin Wang, Xiaoling Wang, 1253 and Liang He. 2024. A safety realignment frame-1254 work via subspace-oriented model fusion for large 1255 language models. arXiv preprint arXiv:2405.09055. 1256

1257

1258

1259

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

- Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jentse Huang, Jiahao Xu, Tian Liang, Pinjia He, and Zhaopeng Tu. 2024. Refuse whenever you feel unsafe: Improving safety in llms via decoupled refusal training. arXiv preprint arXiv:2407.09121.
- Qiusi Zhan, Richard Fang, Rohan Bindu, Akul Gupta, Tatsunori Hashimoto, and Daniel Kang. 2023. Removing rlhf protections in gpt-4 via fine-tuning. ArXiv, abs/2311.05553.
- Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan Cheng, Ziwen Xu, Xin Xu, Jia-Chen Gu, Yong Jiang, Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang, Xiaowei Zhu, Jun Zhou, and Huajun Chen. 2024a. A comprehensive study of knowledge editing for large language models. Preprint, arXiv:2401.01286.
- Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. 2024b. Instruction tuning for large language models: A survey. Preprint, arXiv:2308.10792.
- Weixiang Zhao, Yulin Hu, Zhuojun Li, Yang Deng, Yanyan Zhao, Bing Qin, and Tat-Seng Chua. 2024. Towards comprehensive and efficient post safety alignment of large language models via safety patching. arXiv preprint arXiv:2405.13820.
- Chujie Zheng, Fan Yin, Hao Zhou, Fandong Meng, Jie Zhou, Kai-Wei Chang, Minlie Huang, and Nanyun Peng. 2024. On prompt-driven safeguarding for large language models.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. Preprint, arXiv:2306.05685.
- Minjun Zhu, Linyi Yang, and Yue Zhang. 2024. Personality alignment of large language models. arXiv preprint arXiv:2408.11779.
- Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. 2023a. Representation engineering: A topdown approach to ai transparency. arXiv preprint arXiv:2310.01405.
- Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, J Zico Kolter, Matt Fredrikson, and Dan Hendrycks. 2024. Improving alignment and robustness with circuit breakers. In The Thirty-eighth Annual Conference on Neural Information Processing Systems.

1311

1313

1314

1315

1317

1318

1319 1320

1321

1322

1323

1324

1325

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1339

1340

1341

1342

1343

1344

1345

1346

1347

1349

1350

1352

1353

1354

1355

1357

## Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, and Matt Fredrikson. 2023b. Universal and transferable adversarial attacks on aligned language models. Preprint, arXiv:2307.15043.

# A Robustness of SafetyLock against fine-tunning

We examined the safety directions  $\theta_l^h$  in both the original Llama-3-Instruct 8B model and its fine-1316 tuned variants subjected to different risk levels. Focusing on the most effective attention head (the 26th head in the 31st layer) for clarity, as depicted in Figure 5, we observed distinct clustering of activations corresponding to safe (blue) and unsafe (orange) responses across both original and finetuned models. The black arrows in Figures 5a-d illustrate that the shift from unsafe to safe activations maintains a high degree of similarity and consis-1326 tency, regardless of the fine-tuning risk parameters applied. Additionally, our quantitative analysis using cosine similarity (Figure 5e-g) revealed that the similarity between the original and fine-tuned models remains exceptionally high (above 0.99) across all tested risk levels. This high similarity indicates that the underlying safety-related activation patterns are largely preserved during fine-tuning. Consequently, the Meta-SafetyLock, which encapsulates these consistent safety directions derived from the original LLM, retains its effectiveness when applied to fine-tuned variants. This inherent preservation of safety activation patterns eliminates the 1338 need for recalibration, allowing Meta-SafetyLock to generalize seamlessly across different fine-tuned models.

#### **Consistency of Harmlessness Directions** B in Fine-tuned Models

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

We first extracted activation values from the 31st layer, 26th head of the Llama-3-8B Instruct model, which we identified as the most sensitive to harmlessness through linear regression, achieving the highest binary classification accuracy. We then performed forward computation on a safety dataset, saving the activation values of the last token for both safe and unsafe samples. Using 2D PCA for dimensionality reduction, we visualized the shift

in activation values between safe and unsafe samples by connecting their center points with arrows, illustrating both the direction and magnitude of the shift.

1358

1359

1360

1361

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

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

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

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

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

$$\theta_{\alpha} = \theta + \alpha(\theta' - \theta)$$
(5) 13

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

The interpolation process is implemented as follows:

- 1. We first extract the state dictionaries of both the base model  $(\theta)$  and the fine-tuned model  $(\theta').$
- 2. For each layer, we compute the difference vector:  $d_1 = \theta' - \theta$ .
- 3. We then create new weights for each  $\alpha$  value:  $\theta_{\alpha} = \theta + \alpha d_1.$
- 4. These new weights are used to reconstruct a 1397 new state dictionary, maintaining the origi-1398 nal structure and naming conventions of the model. 1400



Figure 5: Analysis of safety directions at the 31st layer, 26th head for the original and fine-tuned models under different risk levels. (a-d) Activation density distributions. (e-g) Cosine similarity plots.

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

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1425

1426

# C Mathematical Explanation of SafetyLock's Effectiveness in Suppressing Harmful Outputs

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

#### C.1 Activation Space and Safety Directions

1420Let us denote the activations of the original (pre-1421fine-tuned) language model at layer l and head h as1422 $\mathbf{x}_{l,h} \in \mathbb{R}^D$ , where D is the dimensionality of the1423head's output. During inference, these activations1424encode information about the generated tokens.

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

$$\mathcal{X}_{\text{safe}} = \left\{ \mathbf{x}_{l,h}^{\text{safe},i} \right\}_{i=1}^{N_{\text{safe}}}, \quad (6) \quad 142$$

$$\mathcal{X}_{\text{unsafe}} = \left\{ \mathbf{x}_{l,h}^{\text{unsafe},i} \right\}_{i=1}^{N_{\text{unsafe}}}, \quad (7) \quad 1428$$

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

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

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

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

# C.2 Preservation of Safety Directions During Fine-Tuning

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

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

1434

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1429

1430

1431

1432

between the activations of the original and fine-tuned models:

$$\tilde{\mathbf{x}}_{l,h} \approx \mathbf{x}_{l,h} + \Delta \mathbf{x}_{l,h},\tag{9}$$

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

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

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

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

## C.3 Effectiveness of Activation Intervention

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

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

where:

1453

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464

1465

1466

1467

1468

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

- α ∈ ℝ is the scaling factor controlling the intensity of the intervention.
- σ<sub>l,h</sub> ∈ ℝ<sup>D</sup> is the standard deviation vector of activations along each dimension, capturing the typical variability.
- $\odot$  denotes element-wise multiplication.

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

## C.4 Impact on Output Probabilities

1484The language model generates the next token based1485on a probability distribution computed from the1486final activations. Adjusting the activations as in1487Equation equation 11 influences the logits  $\mathbf{z} \in \mathbb{R}^V$ 1488(where V is the vocabulary size) before the softmax1489function:

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

1491

1492

1493

1494

1495

1496

1497

1498

1499

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1512

1513

1514

1515

1516

1517

1518

1519

1520

1521

1522

1523

1524

1525

1526

1527

1528

1529

1530

1532

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

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

## C.5 Suppressing Harmful Outputs

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

$$P(t_{\text{harm}}) = \frac{\exp\left(z_{t_{\text{harm}}}^{\text{intervened}}\right)}{\sum_{i=1}^{V} \exp\left(z_{i}^{\text{intervened}}\right)}.$$
 (13) 150

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

### C.6 Transferability Across Models

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

This property is supported by empirical observations of low Kullback–Leibler (KL) divergence between the activation distributions of the original and fine-tuned models (see Figure 5 in Section A). The minimal divergence indicates that the overall structure of the activation space, especially along dimensions relevant to safety, is preserved during fine-tuning.

### C.7 Conclusion

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

1571

1572

1575

1576

1578

1579

1580

1533

for efficient and transferable safety interventions across different models and fine-tuning scenarios.

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

# D The Risks of Fine-tuning LLMs and Experimental Setup

HEx-PHI (Qi et al., 2024) is based on 11 categories of prohibited use cases merged from Meta's Llama-3 acceptable use policy and OpenAI's usage policies: (1) Illegal Activity, (2) Child Abuse Content, (3) Hate, Harass, Violence, (4) Malware, (5) Physical Harm, (6) Economic Harm, (7) Fraud, Deception, (8) Adult Content, (9) Political Campaigning, (10) Privacy Violation Activity, and (11) Tailored Financial Advice. The dataset includes 30 examples per category, totaling 330 examples. This ensures a comprehensive safety evaluation aligned with industry-standard usage policies.

For Risk-1, we use negative samples from the HH-RLHF preference dataset. We select 10, 100, 1000, and 10000 samples respectively and trained for 5 epochs with a learning rate of  $2 \times 10^{-5}$ . For Risk-2, we use 10 samples from Qi et al. (2024) and trained for 5 epochs with a learning rate of  $2 \times 10^{-5}$ . For Risk-3, we use the first 50,000 samples from the Alpaca dataset (Wang et al., 2023b) and trained for 5 epochs with a learning rate of  $2 \times 10^{-5}$ . We set the last token r = 5.

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

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

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

## **E** Additional Experiments

#### E.1 Analysis of SafetyLock's Intervention

**Distance**  $\alpha$ . Our experimental results, as illustrated in Figure 6, demonstrate the significant influence of SafetyLock's intervention distance ( $\alpha$ ) 1610 on model safety across different model sizes. For 1611 both Llama-3-8B and Llama-3-70B, we observe a 1612 clear U-shaped trend in harmfulness metrics as  $\alpha$ 1613 increases. Initially, as  $\alpha$  rises from 0 to 4, there's 1614 a sharp decrease in harmfulness scores and rates, 1615 as well as the AdvBench ASR. This indicates that 1616 moderate intervention effectively enhances model 1617 safety. However, beyond  $\alpha = 4$ , we see a gradual 1618 increase in these metrics, suggesting that exces-1619 sive intervention may lead to unintended conse-1620 quences, potentially disrupting the model's learned 1621 safety boundaries. Notably, Llama-3-70B exhibits more stability across different  $\alpha$  values compared 1623 to Llama-3-8B, implying that larger models may be 1624 more resilient to intervention adjustments. These 1625 findings underscore the importance of carefully cal-1626 ibrating SafetyLock's intervention parameters to 1627 achieve optimal safety improvements while main-1628

1581 1582 1583

1584

1586

1587

1588

1590

1592

1593

1594

1595

1596

1598

1599

1600

1604 1605

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



Figure 6: Impact of SafetyLock's intervention distance ( $\alpha$ ) on model safety metrics for Llama-3-8B and Llama-3-70B models. The graphs show Harmfulness Average Score, Harmfulness Average Rate, and AdvBench ASR across different  $\alpha$  values. Note that for these experiments, the intervention degree K is set to 24, indicating the number of attention heads influenced by SafetyLock.

taining model performance, with an optimal  $\alpha$  value around 4-6 for both model sizes.

1629

1630

1632

1634

1635

1636

1637

1638

1639

1640

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

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

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

1641Our analysis reveals a nuanced pattern of safety1642improvement across different model scales. For1643Llama-3-8B and Llama-3-70B, we observe a rapid1644enhancement in safety metrics as K increases from16450 to 6, followed by more gradual improvements1646up to K=24. This pattern holds consistent across1647all measured metrics: Harmfulness Average Score,

Harmfulness Average Rate, and AdvBench ASR. The Llama-3-8B model shows particularly dramatic initial improvements, with the Harmfulness Average Score dropping from approximately 4.0 to 1.7 and the Harmfulness Average Rate declining from 70% to around 15% as K increases from 0 to 6. The Llama-3-70B model demonstrates similar trends but with generally lower baseline harmfulness scores, suggesting that larger models might possess inherently stronger safety characteristics. Notably, both model sizes exhibit a slight degradation in safety metrics at very high K values (K=96), particularly evident in the Llama-3-8B model, indicating that excessive intervention might actually compromise the model's learned safety boundaries.

1648

1649

1651

1652

1654

1656

1657

1658

1659

1660

Through these experiments, we've identified a consistent scaling pattern across model sizes: 1B-1664 scale models achieve optimal performance with K = 6-12 heads, 8B-scale models with K = 12-241666 heads, and 70B/123B-scale models with K = 24-481667 heads. This scaling law reveals that the proportion of safety-sensitive attention heads actually decreases as model size increases, with larger models 1670 requiring a smaller relative proportion of heads for effective safety control. The identification of this scaling relationship enables direct determination 1673 of appropriate K values based on model size with-1674 out additional search time, significantly enhanc-1675 ing SafetyLock's deployment efficiency. These findings demonstrate that targeted intervention on 1677



Figure 7: Impact of SafetyLock's intervention degree (K) on model safety metrics for Llama-3-8B and Llama-3-70B models. The graphs illustrate the Harmfulness Average Score, Harmfulness Average Rate, and AdvBench ASR across different K values, ranging from 0 to 96. Lower scores indicate better safety performance. Note the rapid improvement in safety metrics as K increases from 0 to 6, followed by more gradual enhancements up to K=24, with a slight uptick at K=96 for some metrics.

a carefully selected subset of attention heads can achieve substantial safety improvements without requiring extensive architectural modifications, highlighting the efficiency and effectiveness of our approach.

# E.2 Impact of Learning Rate on Safety Degradation

1678

1679

1683

1684

1686

1687

1689

1691

1692

1693

1694

1695

1697

1698

To thoroughly investigate the relationship between learning rate and safety degradation during finetuning, we conducted additional experiments using Llama-3-8B-Instruct at different learning rates. Following the hyperparameter settings from previous work (Qi et al., 2024) (detailed in Appendix G.1), we initially used a learning rate of 2e-5 for our main experiments. However, considering that smaller learning rates (e.g., 1e-6) are commonly used in continued pre-training scenarios to minimize impact on model behaviors, we performed comparative experiments under Risk Level-3 finetuning scenario.

 Table 4: Impact of Learning Rate on Safety Degradation

 and Recovery

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

Results in Table 4 demonstrate that a lower learn-

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

1700

1701

1703

1704

1705

1706

1708

1709

1710

1712

1713

1714

1715

1716

1717

1718

1719

1720

1721

1722

1723

1724

1726

# E.3 Direction Consistency Across Multiple Attention Heads

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

Notably, across all examined heads, we observe consistent directional patterns between unsafe and

1727safe content centers, regardless of the fine-tuning1728condition. This consistency validates our core hy-1729pothesis that safety-related patterns in attention1730heads remain largely preserved during fine-tuning,1731enabling effective deployment of Meta-SafetyLock1732extracted from the base model to various fine-tuned1733variants.

1734

1735

1736

1737

1738

1740

1741

1742

1743

1744

1745

1746

1747

1748

1749

1750

1751

1752

1753

1754

1756

1757

1758

1759

1760

1761

1762

1765

1766

1767

1768

1769

# E.4 Domain-Specific Performance: A Case Study on Mathematical Reasoning

To rigorously evaluate SafetyLock's ability to maintain domain-specific capabilities while ensuring safety, we conducted extensive experiments using the GSM8K dataset, a challenging mathematical reasoning benchmark. We fine-tuned Llama-3-8B-Instruct on GSM8K's training set and evaluated both safety metrics and mathematical performance.

Table 5: Safety and Performance Metrics for GSM8KFine-tuning

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

As shown in Table 5, SafetyLock demonstrates remarkable effectiveness in preserving mathematical reasoning capabilities while enhancing safety measures. The minimal performance drop in GSM8K accuracy (from 85.59% to 84.91%) stands in stark contrast to traditional safety-alignment methods like Model-Edited (DINM), which suffers catastrophic degradation to 5.00% accuracy. Simultaneously, SafetyLock achieves superior safety metrics, reducing AdvBench ASR from 7.23% to 0.19% and improving the HEx-PHI Score from 1.45 to 1.12. These results provide compelling evidence that SafetyLock can successfully maintain domain-specific capabilities while ensuring robust safety guardrails, addressing a critical challenge in deploying safe and effective language models for specialized tasks.

# E.5 Impact of Activation Normalization on SafetyLock

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

Results in Table 6 demonstrate the critical role of  $\sigma_l^h$  in balancing safety and model utility. With-

out normalization, while safety metrics improve 1770 marginally (ASR: 0.0%, HEx-PHI: 1.03), the 1771 model suffers severe performance degradation on 1772 GSM8K (52.24%). Including  $\sigma_l^h$  maintains strong 1773 safety improvements while preserving the model's 1774 mathematical reasoning capabilities (84.91% accu-1775 racy). This suggests that activation-specific scal-1776 ing through  $\sigma_l^h$  is essential for preventing over-1777 aggressive interventions that could compromise 1778 model functionality. These findings validate our 1779 design choice and highlight the importance of care-1780 ful calibration in safety interventions. 1781

1782

1783

1784

1785

1786

1787

1788

1789

1790

1791

1792

1793

1794

1795

1796

1797

1798

1799

1800

1801

1802

1803

1804

1805

1806

1807

1808

1809

1810

1811

1812

1813

1814

1815

1816

1817

1818

1819

1820

## E.6 Comparison with Circuit Breakers

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

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

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

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

The performance disparity may be attributed to

Table 6: Impact of Activation Normalization on Safety and Performance

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

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

Method	Level-1 (ASR/Score)	Level-2 (ASR/Score)	Level-3 (ASR/Score)
Original Fine-tuned	49.24%/4.13	38.46%/3.19	42.88%/3.23
Circuit Breakers	84.62%/3.62	27.12%/2.10	94.04%/3.79
SafetyLock	0.19%/1.36	5.19%/1.07	0.19%/1.04

Table 8: Performance on Circuit Breakers' originalbenchmark scenarios

Method	AutoDAN	PAIR	GCG
Base Model	3.7%	18.7% 7.5%	44.5% 2.5%
SafetyLock	0.0% 0.0%	0.0%	2.3% <b>0.0%</b>

Circuit Breakers' representation remapping strategy being less effective when model safety boundaries have been substantially modified through finetuning. SafetyLock's approach of targeting specific attention heads appears more robust to such modifications while maintaining computational efficiency.

# E.7 Impact of Token Window Size on SafetyLock

1821

1822

1824

1825

1826

1827

1829

1830

1831

1832

1833

1834

1835

1837

1838

1839

1840

1841

1842

1843

1845

The choice of how many final tokens to consider when calculating safety directions represents a crucial design decision in SafetyLock's implementation. While previous works often use the entire hidden state for intervention, we hypothesized that focusing on a smaller window of final tokens might capture safety-relevant patterns more effectively while maintaining computational efficiency.

To determine the optimal token window size, we conducted extensive experiments varying r from 1 to 10 tokens across all three risk levels, as shown in Table 9. Our findings reveal that r = 5 consistently achieves optimal or near-optimal safety performance across all scenarios. While smaller windows (r = 1, 3) can effectively improve safety, they may not capture sufficient context for robust

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

 Performance

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

intervention. Conversely, larger windows (r = 10) show slightly degraded performance, possibly due to including less relevant contextual information. This empirical evidence supports our choice of r =5 as the default parameter, offering the best balance between robust safety improvement and effective intervention across different fine-tuning scenarios.

1846

1847

1848

1851

1852

1853

1854

1855

1856

1857

1858

1860

1861

1862

1863

# E.8 Comparison of safety performance within each category

The radar charts in Figure 9 illustrate SafetyLock's effectiveness across eleven distinct safety attack categories for each risk level and model size. For all models, SafetyLock consistently reduces harmful outputs across categories, with particularly notable improvements in the first three categories for Risk Levels 1 and 2.

# F Recommendations for Deploying SafetyLock

Understanding the diverse landscape of model de-1864ployment scenarios is crucial for effectively imple-1865menting SafetyLock to maintain safety while en-1866

1879

1880

1881

1882

1883

1884

1885

1886

1887

1888

1889

1890 1891

1892

1893

1894

1895

1896

1897

1898

1899

1900 1901

1902

1903 1904

1905

1906

1907

1908

1910

1911

1912

1913 1914

1915

1916

1917

1918

1867

abling customization. The method's effectiveness and implementation strategy vary significantly depending on the model's distribution approach and user priorities, leading to distinct considerations for different deployment contexts.

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

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

To address the fundamental challenge of malicious users with full access to open-source weights, we propose a hybrid deployment strategy that combines transparency with controlled access to safetycritical components. This approach involves opensourcing the majority of model weights while retaining control of a small subset of safety-critical weights using methods like Taylor Unswift (Wang et al., 2024a). By providing efficient access to these controlled weights through a service API and applying SafetyLock during the serving phase, organizations can maintain crucial safety controls while preserving the benefits of open-source accessibility. This balanced solution ensures that users can customize models for their specific needs without easily circumventing safety measures, as the critical safety-related parameters remain protected under controlled access. 1919

1920

1921

For successful implementation, organizations 1922 should establish comprehensive monitoring sys-1923 tems to regularly update safety vectors, implement 1924 automatic safety checks post-fine-tuning, and de-1925 velop clear protocols for handling potential con-1926 flicts between safety measures and legitimate use 1927 cases. Regular assessment and updating of safety 1928 mechanisms ensure that SafetyLock remains ef-1929 fective against evolving harmful behaviors, while 1930 clear documentation and guidelines help users un-1931 derstand the implications and importance of these 1932 safety measures. Through these carefully con-1933 sidered deployment strategies and best practices, 1934 SafetyLock provides a robust framework for main-1935 taining model safety across various deployment 1936 scenarios, acknowledging and addressing the inher-1937 ent challenges in protecting open-source models 1938 while enabling their beneficial applications.



Figure 8: Visualization of activation patterns for multiple attention heads. Each row represents a different attention head position, showing consistent directional patterns across the original model and fine-tuned variants. The black arrows indicate the direction from unsafe to safe content centers, demonstrating remarkable consistency in safety directions despite fine-tuning modifications.



Figure 9: Safety performance comparison for 3 Risk Levels fine-tuned LLMs. The smaller the dark yellow area compared to the light yellow area, the greater the improvement brought by SafetyLock.