# SafeInt: Shielding Large Language Models from Jailbreak Attacks via Safety-Aware Representation Intervention

Anonymous ACL submission

#### Abstract

With the widespread real-world deployment of large language models (LLMs), ensuring their behavior complies with safety standards has become crucial. Jailbreak attacks exploit vulnerabilities in LLMs to induce undesirable behavior, posing a significant threat to LLM safety. Previous defenses often fail to achieve both effectiveness and efficiency simultaneously. Defenses from a representation perspective offer new insights, but existing interventions cannot dynamically adjust representations based on the 011 harmfulness of the queries. To address this limi-012 tation, we propose SafeIntervention (SafeInt), 014 a novel defense method that shields LLMs from jailbreak attacks through safety-aware representation intervention. Built on our analysis of the representations of jailbreak samples, the core idea of SafeInt is to relocate jailbreak-related 019 representations into the rejection region. This is achieved by intervening in the representation distributions of jailbreak samples to align them with those of unsafe samples. We conduct comprehensive experiments covering six jailbreak attacks, two jailbreak datasets, and two utility benchmarks. Experimental results demonstrate that SafeInt outperforms all baselines in defending LLMs against jailbreak attacks while largely maintaining utility. Additionally, we evaluate SafeInt against adaptive attacks and verify its effectiveness in mitigating real-time attacks. WARNING: This paper may contain content that is offensive and harmful.

### 1 Introduction

034

039

042

Large Language Models (LLMs) (OpenAI et al., 2024; Touvron et al., 2023; Grattafiori et al., 2024) have demonstrated remarkable performance across various domains (Zhang et al., 2023; Liu et al., 2023; Wang et al., 2024). With their widespread application in real-world scenarios, LLMs face safety challenges (Ferrara, 2023; Ji et al., 2023). Although efforts (Wang et al., 2023a; Rafailov et al., 2024; Zhou et al., 2023) have been made to align LLMs'

behaviors with human values through carefully designed training strategies, recent studies (Zou et al., 2023b; Li et al., 2024b; Mehrotra et al., 2024; Liu et al., 2024b) reveal that LLMs can still produce undesirable behaviors when subjected to well-crafted jailbreak attacks, such as the biased generation or potentially harmful responses. 043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

Various defense methods have been proposed to address the growing threat of jailbreak attacks. Prompt-based defenses use instructions (Phute et al., 2024; Xie et al., 2023; Zhang et al., 2024b) or context (Zhou et al., 2024a; Wei et al., 2024) to prevent LLMs from generating harmful content. However, prompt-based methods rely on manually crafted secure prompts and possibly lead to excessive self-censorship (Varshney et al., 2024), reducing the helpfulness of LLMs for benign queries. Detection-based defenses compute the perplexity of inputs (Alon and Kamfonas, 2023) or perturb them (Cao et al., 2024) to identify jailbreak prompts. Decoding-based defenses (Xu et al., 2024; Liu et al., 2024a) reconstruct a safer output probability distribution through contrastive decoding. However, these methods often lack effectiveness or require additional inference overhead. We aim to defend LLMs against jailbreak attacks from a representation perspective, which provides a more controllable and efficient approach. Previous studies (Zou et al., 2023a; Rimsky et al., 2024) have shown the effectiveness of intervening representations to steer LLMs' behaviors, but such interventions cannot dynamically adjust representations based on whether a query is harmful. This limitation makes it challenging to leverage representations for mitigating jailbreak attacks.

In this paper, we analyze the representations of jailbreak samples on four LLMs. Our analysis uses a classifier as a proxy to investigate whether jailbreak representations are distinguishable and whether the representation distributions of different jailbreak methods are consistent. We derive two observations. First, in both intermediate and later layers of LLMs, the representations of jailbreak samples can be distinguished from those of safe or unsafe samples. Second, the consistency of the representation distributions across different jailbreak methods is observed in all LLMs, and it is generally more pronounced in the intermediate layers.

086

090

100

101

102

103

104

106

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

127

128

129

130

131

132

133

134

135

Building on these observations, we propose SafeIntervention (SafeInt), a novel defense method that shields LLMs from jailbreak attacks via safety-aware representation intervention. The representations of unsafe samples inherently characterize the rejection region of the LLM. However, jailbreak samples often produce representations that deviate from those of unsafe samples, causing the model to fail to trigger its built-in rejection behavior. The core idea of SafeInt is to relocate jailbreak-related representations into the rejection region, thereby activating the model's native refusal mechanisms. To achieve this, we first project the representations at an intermediate layer into a linear subspace, followed by a parameterized intervention. For jailbreak-related representations, we align their distribution with that of unsafe samples across the subsequent layers. For jailbreakirrelevant representations, we perform representation reconstruction to preserve their original semantics. After training, SafeInt can adaptively intervene in jailbreak-related representations while seamlessly integrating into the LLM inference process.

We conduct a comprehensive evaluation of SafeInt, covering six jailbreak attacks, two jailbreak datasets, and two utility benchmarks. Experimental results show that SafeInt consistently outperforms all baselines in defending against jailbreak attacks. In most cases, SafeInt also maintains the best utility. Additionally, we evaluate SafeInt against adaptive attacks and verify the effectiveness of SafeInt in defending against real-time attacks. In summary, our main contributions are as follows:

- We observe that the representations of jailbreak samples are distinguishable and that the representation distributions of different jailbreak methods exhibit consistency.
- We propose SafeInt, a novel defense method that can adaptively identify and intervene in jailbreak-related representations to shield LLMs from jailbreak attacks.
- Extensive experiments show that SafeInt significantly outperforms all baselines in defend-

ing against jailbreak attacks while largely maintaining utility.

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

### 2 Preliminaries

### 2.1 Representation Intervention

Representation intervention is an effective means of steering LLM behavior. For a given decoderonly transformer model with L layers, we denote the internal representation (or residual stream activation) of the last token at layer l as  $\mathbf{h}^{(l)} \in \mathbb{R}^d$ . A typical form of representation intervention is:

$$\widetilde{\mathbf{h}}^{(l)} = \mathbf{h}^{(l)} \pm \epsilon \cdot \mathbf{v}.$$
 (1)

Here,  $\widetilde{\mathbf{h}}^{(l)}$  is the intervened representation,  $\epsilon \in \mathbb{R}$  represents the intervention strength, and  $\mathbf{v} \in \mathbb{R}^d$  denotes the intervention direction.

### 2.2 Analysis of Jailbreak Sample Representations

Recent works have investigated the representation distributions of unsafe and safe samples within LLMs, utilizing their distributional characteristics to enhance safety or facilitate jailbreaks. In this paper, we analyze the representation distributions of three types of samples after introducing jailbreak samples. We construct three training datasets:  $\mathcal{D}_{iailbreak}$ , which consists of jailbreak instructions generated only using GCG (Zou et al., 2023b) on AdvBench (Zou et al., 2023b);  $\mathcal{D}_{unsafe}$ , which includes harmful instructions extracted from MaliciousInstruct (Huang et al., 2023) and TDC2023 (Mazeika et al., 2023); and  $\mathcal{D}_{safe}$ , which contains harmless instructions sampled from Alpaca (Taori et al., 2023).<sup>1</sup> More details of the datasets are provided in Appendix A.1. We conduct our analysis on four LLMs: Qwen-7B-Chat (Bai et al., 2023), Llama-2-7B-Chat (Touvron et al., 2023), Llama-3-8B-Instruct (Grattafiori et al., 2024), and Vicuna-7B-v1.5 (Chiang et al., 2023). In each layer of the LLM, we train a logistic regression classifier to fit the representations of the three types of samples and report the test accuracy.

*Q1: Are the representations of jailbreak, unsafe, and safe samples distinguishable?* 

In Figure 1, we present the classification accuracy on a test set containing only GCG jailbreak samples. Using the classifier as a proxy for observation, a higher classification accuracy indicates that the representations of the three types of samples are

<sup>&</sup>lt;sup>1</sup>For convenience, we abbreviate these datasets as  $\mathcal{D}_j$ ,  $\mathcal{D}_u$ , and  $\mathcal{D}_s$  in the following.



Figure 1: Test accuracy of classifiers at different layers of LLMs, with the test set containing only GCG jailbreak samples.



Figure 2: Test accuracy of classifiers on a test set containing jailbreak samples from GCG, AutoDAN, and DeepInception.

more distinguishable. For all LLMs, the test accuracy remains above 95% starting from the 10th or 11th layer. This indicates that the representations of the three types of samples become distinguishable from the intermediate layers of LLMs onward.

183

184

185

190

191

192

193

194

195

196

197

203

Q2: Are the representation distributions of samples generated by different jailbreak methods consistent?

We reconstruct a test set where jailbreak samples are composed of three methods: GCG, AutoDAN (Liu et al., 2024b), and DeepInception (Li et al., 2024b). We employ the previously trained classifiers for testing and show the results in Figure 2. Since the classifiers are trained solely on GCG jailbreak samples, a high test accuracy reveals that the representations generated by different jailbreak methods exhibit a unified pattern from the classifier's perspective, indicating consistency.

We observe this consistency across different LLMs. For Qwen, Llama2, and Llama3, the accuracy remains above 90% in most layers. For Vicuna, the accuracy exceeding 90% is primarily observed in the intermediate layers. Although the trend of consistency across layers varies among different LLMs, it is generally more pronounced in the intermediate layers. 204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

**Key Insights and Motivation** Aligned LLMs can reject unsafe samples, and the representations of these samples inherently characterize the model's rejection region. Since the representations of jailbreak samples differ from those of unsafe samples, they need to be relocated into the rejection region. Based on the distinguishability and distributional consistency of jailbreak representations, we aim to intervene in their representations to align their distribution with that of unsafe samples.

### 3 Method

In this section, we describe how SafeInt enhances the safety of LLMs. Figure 3 illustrates the diagram of SafeInt.

### 3.1 Representation Relocation

We achieve representation relocation by a targeted intervention that maps jailbreak-related representations into the rejection region defined by unsafe samples. According to the linear interpretability hypothesis commonly used in existing methods (Zhang et al., 2024a; Li et al., 2024a), deep model embeddings can be linearly transformed to correspond to specific human concepts. Thus, we aim to apply a parameterized intervention within a representation space that corresponds to safety-relevant concepts, minimizing impacts on other capabilities. Inspired by LoReFT (Wu et al., 2024), we project the representations into a linear subspace defined by a low-rank projection matrix. Assuming the intervention is applied at layer  $\mathcal{I}$ , it can be parameterized as follows:

$$\widetilde{\mathbf{h}}^{(\mathcal{I})} = \mathbf{h}^{(\mathcal{I})} + \mathbf{U}^{\top} \left( f_{\boldsymbol{\theta}}(\mathbf{h}^{(\mathcal{I})}) - \mathbf{U}\mathbf{h}^{(\mathcal{I})} \right). \quad (2)$$

The matrix  $\mathbf{U} \in \mathbb{R}^{r \times d}$  has orthonormal rows, where *r* denotes the rank of the subspace. The function  $f_{\theta}$  is a linear relocation mapping defined as  $f_{\theta} : \mathbb{R}^d \to \mathbb{R}^r$ .

Then, we define the objectives for learning the intervention. Broadly, our objectives are twofold: a **safety objective** and a **utility objective**. The safety objective guarantees the intervention to help the LLM reject jailbreak and harmful instructions. The utility objective ensures that the intervention does not degrade the response quality for harmless instructions.



Figure 3: The schematic of SafeInt. We apply the intervention (illustrated by the blue shield) at a specific layer and perform alignment in the subsequent layers. The distribution of jailbreak sample representations is adjusted to align with that of unsafe samples while minimizing shifts in the representations of safe and unsafe samples. With the original representation distribution, the LLM is successfully jailbroken and generates harmful content. After alignment, the LLM safely rejects the jailbreak instruction.

#### 3.2 Representation Alignment

252

255

256

260

261

265

267

272

274

We use the classifier as a proxy to assess whether the distributions of jailbreak samples and unsafe samples are consistent in the representation space. From the perspective of the classifier, the alignment is achieved when the classification results for jailbreak and unsafe sample representations are consistent. Specifically, for jailbreak samples, we intervene on their representations to maximize the probability of being classified as unsafe. For unsafe sample representations, they should still be classified as unsafe with a high probability.

We denote the sets of original representations of  $\mathcal{D}_j$ ,  $\mathcal{D}_u$ , and  $\mathcal{D}_s$  as  $H_j$ ,  $H_u$ , and  $H_s$ , respectively. The sets of intervened representations are denoted as  $\widetilde{H}_j$ ,  $\widetilde{H}_u$ , and  $\widetilde{H}_s$ . Let  $\mathbb{L}^a$  be the set of layers to be aligned, with  $\min(\mathbb{L}^a) > \mathcal{I}$ . After applying the intervention at the layer  $\mathcal{I}$ , the updated representation is propagated to the subsequent layers. At layer  $l \in \mathbb{L}^a$ , we extract the latest representation  $\widetilde{\mathbf{h}}^{(l)}$  and compute the following:

$$\mathcal{L}_{cls}^{(l)} = -\frac{1}{|\widetilde{H}_{j}^{(l)}|} \sum_{\widetilde{\mathbf{h}}_{j}^{(l)} \in \widetilde{H}_{j}^{(l)}} \log P_{\mathbf{u}}(\widetilde{\mathbf{h}}_{j}^{(l)}) -\frac{1}{|\widetilde{H}_{\mathbf{u}}^{(l)}|} \sum_{\widetilde{\mathbf{h}}_{\mathbf{u}}^{(l)} \in \widetilde{H}_{\mathbf{u}}^{(l)}} \log P_{\mathbf{u}}(\widetilde{\mathbf{h}}_{\mathbf{u}}^{(l)}), \quad (3)$$

where  $P_{\rm u}$  represents the probability that classifier *P* classifies a representation as unsafe. **Contrastive Learning** Although we align the representations of jailbreak samples with unsafe samples by a classifier, the limited training data may prevent the classifier's decision boundary from accurately capturing the discriminative boundary within the LLM. To enhance the alignment, we use contrastive learning as a complementary task. For a given representation  $\mathbf{q}$ , contrastive learning maximizes the similarity between  $\mathbf{q}$  and the set of positive samples  $K^+$  while minimizing the similarity between  $\mathbf{q}$  and the set of negative samples  $K^-$ , with the objective formulated as follows:

277

278

281

287

288

289

290

294

297

298

300

302

$$CT(\mathbf{q}, K^+, K^-) = -\log \frac{\exp(\operatorname{sim}(\mathbf{q}, \mathbf{k}^+)/\tau)}{\sum_{\mathbf{k} \in (K^+, K^-)} \exp(\operatorname{sim}(\mathbf{q}, \mathbf{k})/\tau)}, \quad (4)$$

where  $\mathbf{k}^+ \in K^+$ ,  $sim(\cdot, \cdot)$  represents cosine similarity, and the temperature is set to  $\tau = 0.1$ .

Specifically, the intervened representations of jailbreak samples should be as close as possible to those of unsafe samples while being pushed away from their original representations and those of safe samples. Accordingly, for  $\widetilde{\mathbf{h}}_{j}^{(l)} \in \widetilde{H}_{j}^{(l)}$ , the contrastive loss is calculated as:

$$\mathcal{L}_{ct}^{(l)} = \text{CT}(\widetilde{\mathbf{h}}_{j}^{(l)}, H_{u}^{(l)}, (H_{j}^{(l)} \cup H_{s}^{(l)})).$$
(5)

### 3.3 Representation Reconstruction

To prevent excessive intervention from distorting the LLM's internal representations, we introduce a

390

391

392

394

395

396

397

399

reconstruction loss to constrain jailbreak-irrelevant representations from changing. Specifically, we encourage the representations of safe and unsafe samples after intervention to remain close to their original states. This ensures that the intervention primarily affects jailbreak-related representations without causing unnecessary shifts in the model's overall representation structure. The loss is formulated as follows:

$$\mathcal{L}_{recon} = \text{MSE}(H_s, H_s) + \text{MSE}(H_u, H_u), \quad (6)$$

where MSE refers to the mean squared error loss.

Considering both the alignment and reconstruction objectives, our final loss is calculated as follows:

$$\mathcal{L}_{total} = \alpha \sum_{l \in \mathbb{L}^a} (\mathcal{L}_{cls}^{(l)} + \mathcal{L}_{ct}^{(l)}) + \beta \mathcal{L}_{recon}.$$
 (7)

Through the two hyperparameters  $\alpha$  and  $\beta$ , we achieve a balance between effective alignment and model stability.

### 4 Experiments

313

314

315

317

318

319

320

324

326

327

328

330

334

335

#### 4.1 Experimental Setup

**Models and Datasets** We primarily evaluate SafeInt on two open-source LLMs: Llama2-7bchat and Vicuna-7b-v1.5. Additionally, we assess its scalability by applying it to a heterogeneous LLM, with results reported in Appendix B.1. For evaluation, we randomly sample 50 instructions from AdvBench (Zou et al., 2023b) as the test set, ensuring no overlap with the training set  $D_j$ . Moreover, to demonstrate that SafeInt is data-agnostic, we construct an out-of-distribution test set consisting of 50 instructions randomly sampled from JailbreakBench (Chao et al., 2024a). Following Xu et al. (2024), we use MT-Bench (Zheng et al., 2023) and Just-Eval (Lin et al., 2023) to evaluate the utility of intervened LLMs.

Jailbreak Attacks Multiple representative jailbreak attacks are employed in our evaluation. These include optimization-based attacks: GCG and AutoDAN, LLM-generated attacks: PAIR (Chao et al., 2024b) and TAP (Mehrotra et al., 2024), and scenario-based attacks: DeepInception. We also consider multilingual mismatch generalization attacks (MG) (Yong et al., 2024), where each instruction in the test set is translated into one of six non-English languages to perform the attacks.

**Baselines** We compare SafeInt with six stateof-the-art defense approaches: PPL (Alon and Kamfonas, 2023), Paraphrase (Jain et al., 2023), Self-Examination (Phute et al., 2024), ICD (Wei et al., 2024), Self-Reminder (Xie et al., 2023), and SafeDecoding (Xu et al., 2024). In addition, we include comparisons with two other baselines that leverage representations for defense, with results presented in Appendix B.2.

**Evaluation Metrics** We use two types of Attack Success Rate (**ASR**) to evaluate defense effectiveness: **ASR-keyword**, which matches predefined refusal keywords, and **ASR-GPT**, which leverages GPT-40-mini to assess whether the LLM generates harmful content relevant to the malicious instruction. Lower ASR values indicate better defense performance. For MT-Bench and Just-Eval, we adopt GPT-based scoring, where Just-Eval evaluates five aspects: helpfulness, clarity, factuality, depth, and engagement.

**Implementation Details** Previous classification results indicate that in the intermediate layers, the representations of various jailbreak samples are relatively consistent and highly distinguishable. To avoid time-consuming searches, we directly select layer  $\mathcal{I} = 12$  as the intervention layer. Since Vicuna lacks harmless alignment, it exhibits weaker safety. Accordingly, we set the second half of the layers as the alignment layers. In contrast, for models like Llama2 that have undergone safety alignment, aligning only the final layer is sufficient. Additional settings and discussions are provided in Appendix A.4 to A.6.

#### 4.2 Main Results

Table 1 presents the ASR results of SafeInt and various baselines on AdvBench. For both Vicuna and Llama2, SafeInt achieves the best performance, reducing ASR to the lowest level among all defense methods under different attacks. Although our training process only utilizes jailbreak samples constructed with GCG, SafeInt effectively defends against other attack strategies, such as PAIR and TAP, which generate adversarial prompts using the LLM. This highlights the generalization capability of our defense, validating our previous observation. Moreover, even against MG attacks, SafeInt significantly lowers ASR, showing that it can generalize to different languages.

Table 2 reports results on another out-ofdistribution test set, JailbreakBench. SafeInt continues to outperform all baselines across different models and attack strategies. This demonstrates its

Madal	Defense	Benchmark↓			Jailbreak A	ttacks↓		
WIGUEI	Derense	AdvBench	GCG	AutoDAN	DeepInception	PAIR	TAP	MG
	No Defense	8% (4%)	90% (92%)	82% (88%)	64% (100%)	54% (60%)	84% (80%)	30% (66%)
	PPL	8% (4%)	26% (30%)	72% (68%)	64% (100%)	52% (58%)	84% (82%)	28% (62%)
	Paraphrase	6% (6%)	18% (20%)	34% (52%)	38% (96%)	36% (38%)	42% (52%)	10% (32%)
V	Self-Examination	2% (0%)	12% (16%)	18% (22%)	34% (74%)	8% (14%)	34% (30%)	6% (34%)
vicuna	ICD	0% (0%)	14% (14%)	40% (36%)	64% (96%)	24% (34%)	44% (44%)	6% (34%)
	Self-Reminder	0% (0%)	4% (6%)	8% (6%)	46% (100%)	26% (32%)	38% (40%)	16% (50%)
	SafeDecoding	0% (0%)	2% (2%)	10% (4%)	0% (0%)	4% (6%)	12% (12%)	12% (40%)
	SafeInt (Ours)	0% (0%)	0% (0%)	2% (2%)	0% (0%)	2% (6%)	8% (10%)	4% (8%)
	No Defense	0% (0%)	30% (32%)	34% (44%)	0% (0%)	2% (10%)	10% (10%)	0% (6%)
	PPL	0% (0%)	0% (2%)	2% (8%)	0% (0%)	2% (8%)	10% (10%)	0% (4%)
	Paraphrase	0% (10%)	0% (22%)	6% (26%)	0% (0%)	2% (30%)	2% (30%)	0% (16%)
T. 1	Self-Examination	0% (0%)	0% (4%)	2% (6%)	0% (0%)	2% (4%)	2% (4%)	0% (0%)
Llama2	ICD	0% (0%)	0% (0%)	0% (0%)	0% (0%)	0% (0%)	0% (0%)	0% (0%)
	Self-Reminder	0% (0%)	0% (2%)	0% (0%)	0% (0%)	2% (4%)	0% (2%)	0% (6%)
	SafeDecoding	0% (0%)	0% (2%)	0% (4%)	0% (0%)	0% (6%)	0% (0%)	0% (0%)
	SafeInt (Ours)	0% (0%)	0% (0%)	0% (0%)	0% (0%)	0% (4%)	0% (0%)	0% (0%)

Table 1: ASR-GPT (outer) and ASR-keyword (in parentheses) for different defense methods on AdvBench. The best results are in **bold**. SafeInt outperforms all baselines across various attacks.

M- 1-1	Defense	Benchmark↓			Jailbreak A	ttacks↓		
Model	Defense	JailbreakBench	GCG	AutoDAN	DeepInception	PAIR	TAP	MG
	No Defense	6% (10%)	74% (96%)	76% (98%)	54% (100%)	42% (46%)	66% (68%)	30% (76%)
	PPL	6% (10%)	20% (30%)	48% (62%)	46% (100%)	38% (48%)	66% (68%)	26% (66%)
	Paraphrase	6% (20%)	18% (40%)	22% (60%)	28% (98%)	16% (36%)	32% (42%)	10% (40%)
Vicuna	Self-Examination	0% (4%)	8% (28%)	26% (48%)	28% (74%)	10% (14%)	28% (30%)	8% (50%)
	ICD	0% (0%)	8% (14%)	42% (42%)	54% (94%)	16% (28%)	32% (42%)	22% (52%)
	Self-Reminder	0% (2%)	4% (4%)	4% (6%)	36% (100%)	14% (20%)	26% (30%)	30% (62%)
	SafeDecoding	0% (0%)	0% (0%)	18% (18%)	0% (0%)	10% (14%)	10% (12%)	14% (36%)
	SafeInt (Ours)	0% (0%)	0% (0%)	4% (6%)	0% (0%)	2% (12%)	4% (12%)	8% (24%)
	No Defense	0% (0%)	24% (34%)	30% (42%)	2% (2%)	0% (6%)	6% (6%)	0% (4%)
	PPL	0% (0%)	2% (2%)	6% (6%)	2% (2%)	0% (4%)	6% (6%)	0% (2%)
	Paraphrase	0% (6%)	0% (28%)	0% (22%)	2% (2%)	0% (24%)	4% (28%)	0% (10%)
T 1	Self-Examination	0% (0%)	2% (2%)	6% (6%)	0% (0%)	0% (4%)	2% (2%)	0% (2%)
Liamaz	ICD	0% (0%)	0% (0%)	0%~(0%)	0% (0%)	0% (0%)	0% (0%)	0% (0%)
	Self-Reminder	0% (0%)	0% (0%)	0%~(0%)	0% (0%)	0% (0%)	0% (0%)	0% (10%)
	SafeDecoding	0% (0%)	0% (0%)	0%~(0%)	0% (0%)	0% (2%)	0% (6%)	0% (0%)
	SafeInt (Ours)	0% (0%)	0% (0%)	0% (0%)	0% (0%)	0% (6%)	0% (0%)	0% (0%)

Table 2: ASR-GPT (outer) and ASR-keyword (in parentheses) on JailbreakBench. The best results are in **bold**. SafeInt consistently achieves the best performance.

robustness to unseen data.

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

While delivering strong defense performance, SafeInt largely preserves the utility of LLMs. As shown in Table 3, SafeInt achieves almost identical scores to the non-defended model in Llama2, whereas ICD and Self-Examination severely degrade utility. For Vicuna, SafeInt results in only a 2% decrease in MT-Bench and a 1% decrease in Just-Eval compared to the non-defended model. In contrast, SafeDecoding leads to 7% drops in both benchmarks. See Appendix C.2 for representative examples.

Since our intervention essentially involves an incremental computation, it can be integrated directly into the forward propagation of the model.

Thus, SafeInt can be seamlessly embedded into the inference process of LLMs. Unlike SafeDecoding, which requires an additional expert model for contrastive decoding, SafeInt introduces virtually no extra computational overhead. A detailed efficiency analysis is provided in Appendix B.3.

415

416

417

418

419

420

421

422

423

424

425

426

427

428

### 4.3 Adaptive Attack

We also consider a scenario where the attacker knows SafeInt and has access to the LLM deployed with it. This means the attacker can dynamically adjust their attack strategies based on the latest defended LLM to bypass the defense mechanisms more effectively. Table 4 presents the experimental results in this scenario, showing that SafeInt main-

Madal	Defense	MT Danah A	Just-Eval ↑						
Model	Defense		Helpfulness	Clear	Factual	Deep	Engaging	Average	
	No Defense	5.21	4.44	4.66	4.38	3.60	3.49	4.11	
	Self-Examination	<u>5.03</u>	4.40	4.65	4.34	3.56	3.47	4.08	
Vicuna	ICD	4.86	4.34	4.61	4.34	3.40	3.32	4.00	
	SafeDecoding	4.84	3.92	4.45	4.19	3.24	3.25	3.81	
	SafeInt (Ours)	5.09	4.40	4.64	4.35	3.49	3.41	<u>4.06</u>	
	No Defense	5.80	4.65	4.78	4.50	4.19	3.90	4.40	
	Self-Examination	1.61	3.21	3.67	3.47	2.92	2.68	3.19	
Llama2	ICD	2.91	3.44	4.08	3.96	3.25	3.24	3.59	
	SafeDecoding	<u>5.68</u>	4.53	4.73	4.42	4.05	3.83	<u>4.31</u>	
	SafeInt (Ours)	5.82	4.62	4.76	4.47	4.13	3.89	4.37	

Table 3: Utility evaluation scores of SafeInt and baselines. The highest and second-highest scores obtained by defense methods are in **bold** and underlined, respectively. SafeInt maintains the best utility in most cases.

Jailbreak Attacks	AdvBench	JailbreakBench	Mathada	Jai	lbreak Attack	MT Danah A	
Adaptive_GCG	ive $GCG = 0\% (0\%) = 0\% (6\%)$	0% (6%)	Wiethous	GCG	AutoDAN	PAIR	MI-Bellen
Adaptive-AutoDAN	0%(0%) 0%(0%)	6% (8%)	No Defense	90%	82%	54%	5.21
			SafeInt	0%	2%	2%	5.09

Table 4: Experimental results of defending against adaptive attacks on Vicuna, with evaluation metrics ASR-GPT and ASR-keyword (in parentheses). 'Adaptive-GCG' and 'Adaptive-AutoDAN' refer to GCG and AutoDAN attacks that are optimized in real-time based on the LLM deployed with SafeInt.

tains strong defensive performance. After SafeInt is deployed, even if GCG and AutoDAN optimize their adversarial prompts in real time, it is difficult to generate threatening attacks.

#### 5 Analyses

#### Ablation Studies 5.1

We conduct ablation studies on the introduced contrastive loss and reconstruction loss to verify their effectiveness. As shown in Table 5, removing contrastive loss increases ASR, indicating its crucial role in enhancing defense performance. Incorporating contrastive loss leads to a decrease in MT-Bench scores, which may be attributed to its impact on the overall representation structure when pulling together or pushing apart local representations. When the reconstruction loss is omitted, representations are more susceptible to excessive intervention, resulting in both defense failure and a significant decline in response quality.

#### Hyperparameter Analysis 5.2

Intervention Layer Choice To understand how the choice of the intervention layer impacts our defense effectiveness, we conduct an analysis. Fig-

Mathada	Jai	lbreak Attack	MT Banch $\uparrow$	
wiethous	GCG	AutoDAN	PAIR	MI-Bellell
No Defense	90%	82%	54%	5.21
SafeInt	0%	2%	2%	5.09
w/o $\mathcal{L}_{ct}$	2%	8%	6%	5.22
w/o $\mathcal{L}_{recon}$	2%	12%	8%	4.09

Table 5: Ablation results of our method on AdvBench and MT-Bench, using ASR-GPT as the metric for Jailbreak Attacks. 'w/o  $\mathcal{L}_{ct}$ ' and 'w/o  $\mathcal{L}_{recon}$ ' denote the removal of contrastive loss and reconstruction loss, respectively.

ure 4(a) displays the ASR-keyword when the intervention layer is set between layers 10 and 20. We observe that intervening in the intermediate layers generally yields better results than intervening in the later layers, which may suggest that these intermediate layers play a more crucial role in jailbreak mechanisms. Notably, when the intervention is applied at layer 13, ASR reaches its lowest point. This finding aligns with our observation in Figure 2, where Vicuna exhibits the highest jailbreak representation consistency at layer 13.

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

Alignment Layer Range We fix the endpoint of the alignment layer range at the final layer and modify the starting point to control its span. In Figure 4(b), we illustrate the results when setting the starting point between layers 15 and 25. We observe that as the starting point shifts to later layers, the defense effectiveness weakens. This may be attributed to the reduced number of aligned layers being insufficient to correct the attack. Overall, while adjusting the alignment layer range impacts defense performance, the effect is not drastic, indicating that our method exhibits a certain degree of

429

430

431

- 442 443
- 444 445

446

447

448

449

450



Figure 4: Analysis of the intervention layer and alignment layer range.

robustness to this hyperparameter.

### 6 Related Work

475

476

477

478

479

480

481

482

483

485

486

487

488

489

490

491

492

493

Jailbreak Attacks Jailbreak attacks aim to bypass alignment or safeguards, forcing LLMs to generate inappropriate content. Early jailbreak attacks (Wei et al., 2023; Yong et al., 2024; Yuan et al., 2024) rely on manually crafted adversarial prompts, which primarily exploit objective competition and mismatched generalization to achieve jailbreaks. Subsequent optimization-based attacks (Zou et al., 2023b; Liu et al., 2024b; Paulus et al., 2024) introduce automated adversarial prompt optimization by leveraging the internal states of LLMs, significantly improving both the success rate and efficiency of jailbreaks. Recent jailbreak attacks (Chao et al., 2024b; Mehrotra et al., 2024; Ding et al., 2024) iteratively rewrite and refine adversarial prompts using one or multiple LLMs, further exposing security vulnerabilities in LLMs.

Jailbreak Defenses To address the challenges 494 posed by jailbreak attacks, numerous defense meth-495 ods have been proposed (Robey et al., 2024; Kumar 496 et al., 2025). Detection-based approaches identify 497 498 adversarial prompts by computing perplexity (Alon and Kamfonas, 2023) or randomly deleting parts 499 of the input (Cao et al., 2024). Some methods prompt the LLM to perform self-checking through instructions (Phute et al., 2024; Xie et al., 2023; 502

Zhang et al., 2024b) or context (Zhou et al., 2024a). Decoding-based defenses (Xu et al., 2024; Liu et al., 2024a) focus on analyzing decoding probabilities under different conditions and formulating decoding strategies to ensure safer outputs. Additionally, certain approaches (Zhao et al., 2024) edit specific model parameters to make LLMs forget harmful knowledge. A more controllable and efficient class of defenses (Li et al., 2025; Shen et al., 2025) involves manipulating representations to mitigate jailbreak attacks without modifying model parameters or adding decoding overhead. 503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

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

**Representation Engineering for Safety** Many studies have employed representation engineering techniques (Zou et al., 2023a) to investigate or enhance the safety of LLMs. Zhou et al. (2024b) and Arditi et al. (2024) analyze the mechanisms of jailbreak and refusal from a representation perspective, respectively. Li et al. (2025) improve the robustness of LLMs by strengthening the safety patterns they recognize. Zheng et al. (2024) introduce a learnable safety prompt that aims to increase the separation between harmful and harmless query representations along the refusal direction. Shen et al. (2025) add a difference vector to query representations to guide the LLM toward rejecting malicious instructions, while Gao et al. (2024) mitigate jailbreak attacks by constraining activations within a safe boundary. A major drawback of these two approaches is that their interventions cannot be automatically optimized. This means that when the intervention is applied to all query representations, the choice of intervention strength becomes highly sensitive. In contrast, our method adopts a parameterized intervention, which adaptively identifies and adjusts jailbreak-related representations regardless of manually tuning the intervention strength.

### 7 Conclusion

This paper first analyzes the representations of jailbreak samples on four LLMs and makes key observations. Building on these observations, we propose SafeIntervention (SafeInt), a novel method that defends LLMs against jailbreak attacks via representation intervention. SafeInt can adaptively identify and intervene in jailbreak-related representations while seamlessly integrating into the LLM inference process. Comprehensive experimental results show that our proposed SafeInt outperforms all baselines in defending against jailbreak attacks. In most cases, SafeInt also achieves the best utility.

562

564

565

567

568

571

578

583

584

585

586

587

588

590

594

596

### Limitations

We discuss the limitations of our work. We make a preliminary observation that SafeInt can transfer 555 to different but homologous LLMs. We speculate 556 that these homologous LLMs may share similar jailbreak representation structures. However, we have not conducted an in-depth exploration of the 559 transferability of SafeInt.

### References

- Guillaume Alain and Yoshua Bengio. 2018. Understanding intermediate layers using linear classifier probes. Preprint, arXiv:1610.01644.
- Gabriel Alon and Michael Kamfonas. 2023. Detecting language model attacks with perplexity. Preprint, arXiv:2308.14132.
- Andy Arditi, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. 2024. Refusal in language models is mediated by a single direction. Preprint, arXiv:2406.11717.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, and Xiaodong Deng. 2023. Qwen technical report. Preprint, arXiv:2309.16609.
- Bochuan Cao, Yuanpu Cao, Lu Lin, and Jinghui Chen. 2024. Defending against alignment-breaking attacks via robustly aligned LLM. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 10542-10560, Bangkok, Thailand. Association for Computational Linguistics.
- Patrick Chao, Edoardo Debenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash Sehwag, Edgar Dobriban, Nicolas Flammarion, George J. Pappas, Florian Tramèr, Hamed Hassani, and Eric Wong. 2024a. Jailbreakbench: An open robustness benchmark for jailbreaking large language models. In NeurIPS Datasets and Benchmarks Track.
- Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric Wong. 2024b. Jailbreaking black box large language models in twenty queries. Preprint, arXiv:2310.08419.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An opensource chatbot impressing gpt-4 with 90%\* chatgpt quality.
- Peng Ding, Jun Kuang, Dan Ma, Xuezhi Cao, Yunsen Xian, Jiajun Chen, and Shujian Huang. 2024. A wolf in sheep's clothing: Generalized nested jailbreak prompts can fool large language models easily. Preprint, arXiv:2311.08268.

Emilio Ferrara. 2023. Should chatgpt be biased? challenges and risks of bias in large language models. Preprint, arXiv:2304.03738.

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

- Lang Gao, Xiangliang Zhang, Preslav Nakov, and Xiuying Chen. 2024. Shaping the safety boundaries: Understanding and defending against jailbreaks in large language models. Preprint, arXiv:2412.17034.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, and Alan Schelten, 2024. The llama 3 herd of models. Preprint, arXiv:2407.21783.
- Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. 2023. Catastrophic jailbreak of open-source llms via exploiting generation. Preprint, arXiv:2310.06987.
- 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.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. ACM Comput. Surv., 55(12).
- Aounon Kumar, Chirag Agarwal, Suraj Srinivas, Aaron Jiaxun Li, Soheil Feizi, and Himabindu Lakkaraju. 2025. Certifying llm safety against adversarial prompting. Preprint, arXiv:2309.02705.
- Meng Li, Haoran Jin, Ruixuan Huang, Zhihao Xu, Defu Lian, Zijia Lin, Di Zhang, and Xiting Wang. Evaluating readability and faithfulness 2024a. of concept-based explanations. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 607-625, Miami, Florida, USA. Association for Computational Linguistics.
- Tianlong Li, Zhenghua Wang, Wenhao Liu, Muling Wu, Shihan Dou, Changze Lv, Xiaohua Wang, Xiaoqing Zheng, and Xuanjing Huang. 2025. Revisiting jailbreaking for large language models: A representation engineering perspective. In Proceedings of the 31st International Conference on Computational Linguistics, pages 3158-3178, Abu Dhabi, UAE. Association for Computational Linguistics.
- 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.
- Bill Yuchen Lin, Abhilasha Ravichander, Ximing Lu, Nouha Dziri, Melanie Sclar, Khyathi Chandu, Chandra Bhagavatula, and Yejin Choi. 2023. The unlocking spell on base llms: Rethinking alignment via in-context learning. Preprint, arXiv:2312.01552.

- Quan Liu, Zhenhong Zhou, Longzhu He, Yi Liu, Wei Zhang, and Sen Su. 2024a. Alignment-enhanced decoding: Defending jailbreaks via token-level adaptive refining of probability distributions. In <u>Proceedings</u> of the 2024 Conference on Empirical Methods in <u>Natural Language Processing</u>, pages 2802–2816, Miami, Florida, USA. Association for Computational Linguistics.
  - Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. 2024b. Autodan: Generating stealthy jailbreak prompts on aligned large language models. In <u>The Twelfth International Conference on Learning</u> Representations.
  - Yixin Liu, Avi Singh, C. Daniel Freeman, John D. Co-Reyes, and Peter J. Liu. 2023. Improving large language model fine-tuning for solving math problems. Preprint, arXiv:2310.10047.

678

679

701

706

707

710

711

712

713

714

715

716

- Mantas Mazeika, Andy Zou, Norman Mu, Long Phan, Zifan Wang, Chunru Yu, Adam Khoja, Fengqing Jiang, Aidan O'Gara, Ellie Sakhaee, Zhen Xiang, Arezoo Rajabi, Dan Hendrycks, Radha Poovendran, Bo Li, and David Forsyth. 2023. Tdc 2023 (Ilm edition): The trojan detection challenge. In <u>NeurIPS</u> Competition Track.
- Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron Singer, and Amin Karbasi. 2024. Tree of attacks: Jailbreaking black-box llms automatically. <u>Preprint</u>, arXiv:2312.02119.
- Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Pete Walsh, Oyvind Tafjord, Nathan Lambert, Yuling Gu, Shane Arora, Akshita Bhagia, Dustin Schwenk, David Wadden, Alexander Wettig, Binyuan Hui, Tim Dettmers, Douwe Kiela, Ali Farhadi, Noah A. Smith, Pang Wei Koh, Amanpreet Singh, and Hannaneh Hajishirzi. 2025. Olmoe: Open mixture-of-experts language models. <u>Preprint</u>, arXiv:2409.02060.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, and Florencia Leoni Aleman. 2024. Gpt-4 technical report. <u>Preprint</u>, arXiv:2303.08774.
- Anselm Paulus, Arman Zharmagambetov, Chuan Guo, Brandon Amos, and Yuandong Tian. 2024. Advprompter: Fast adaptive adversarial prompting for llms. Preprint, arXiv:2404.16873.
- Mansi Phute, Alec Helbling, Matthew Hull, ShengYun Peng, Sebastian Szyller, Cory Cornelius, and Duen Horng Chau. 2024. Llm self defense: By self examination, llms know they are being tricked. Preprint, arXiv:2308.07308.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. <u>Preprint</u>, arXiv:2305.18290.

Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. 2024. Steering llama 2 via contrastive activation addition. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 15504–15522, Bangkok, Thailand. Association for Computational Linguistics. 717

718

719

721

723

724

725

726

727

728

729

730

731

732

733

735

736

737

738

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

760

761

762

763

764

765

766

767

768

769

770

771

773

- Alexander Robey, Eric Wong, Hamed Hassani, and George J. Pappas. 2024. Smoothllm: Defending large language models against jailbreaking attacks. Preprint, arXiv:2310.03684.
- Guobin Shen, Dongcheng Zhao, Yiting Dong, Xiang He, and Yi Zeng. 2025. Jailbreak antidote: Runtime safety-utility balance via sparse representation adjustment in large language models. <u>Preprint</u>, arXiv:2410.02298.
- Oscar Skean, Md Rifat Arefin, Yann LeCun, and Ravid Shwartz-Ziv. 2024. Does representation matter? exploring intermediate layers in large language models. Preprint, arXiv:2412.09563.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https:// github.com/tatsu-lab/stanford\_alpaca.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, and Guillem Cucurull. 2023. Llama 2: Open foundation and fine-tuned chat models. Preprint, arXiv:2307.09288.
- Neeraj Varshney, Pavel Dolin, Agastya Seth, and Chitta Baral. 2024. The art of defending: A systematic evaluation and analysis of LLM defense strategies on safety and over-defensiveness. In Findings of the Association for Computational Linguistics: ACL 2024, pages 13111–13128, Bangkok, Thailand. Association for Computational Linguistics.
- Noah Wang, Z.y. Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Jian Yang, Man Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu, Wenhao Huang, Jie Fu, and Junran Peng. 2024. RoleLLM: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. In <u>Findings of the Association for Computational Linguistics: ACL</u> 2024, pages 14743–14777, Bangkok, Thailand. Association for Computational Linguistics.
- Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang, Xin Jiang, and Qun Liu. 2023a. Aligning large language models with human: A survey. <u>Preprint</u>, arXiv:2307.12966.
- Yuxia Wang, Haonan Li, Xudong Han, Preslav Nakov, and Timothy Baldwin. 2023b. Do-not-answer: A dataset for evaluating safeguards in llms. <u>Preprint</u>, arXiv:2308.13387.

866

867

868

869

870

871

872

873

874

832

833

Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2023. Jailbroken: How does llm safety training fail? <u>Preprint</u>, arXiv:2307.02483.

775

778

781

784

787

789

790

791

795

797

802

803

804

805

808

810

811

814

815

816

817

818 819

825

- Zeming Wei, Yifei Wang, Ang Li, Yichuan Mo, and Yisen Wang. 2024. Jailbreak and guard aligned language models with only few in-context demonstrations. <u>Preprint</u>, arXiv:2310.06387.
- Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D. Manning, and Christopher Potts. 2024. Reft: Representation finetuning for language models. <u>Preprint</u>, arXiv:2404.03592.
  - 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. <u>Nature Machine Intelligence</u>, 5(12):1486–1496.
  - Zhangchen Xu, Fengqing Jiang, Luyao Niu, Jinyuan Jia, Bill Yuchen Lin, and Radha Poovendran. 2024. SafeDecoding: Defending against jailbreak attacks via safety-aware decoding. In <u>Proceedings</u> of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long <u>Papers</u>), pages 5587–5605, Bangkok, Thailand. Association for Computational Linguistics.
  - Zheng-Xin Yong, Cristina Menghini, and Stephen H. Bach. 2024. Low-resource languages jailbreak gpt-4. Preprint, arXiv:2310.02446.
  - Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen tse Huang, Pinjia He, Shuming Shi, and Zhaopeng Tu. 2024. Gpt-4 is too smart to be safe: Stealthy chat with llms via cipher. Preprint, arXiv:2308.06463.
  - Hanyu Zhang, Xiting Wang, Xiang Ao, and Qing He. 2024a. Distillation with explanations from large language models. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 5018–5028, Torino, Italia. ELRA and ICCL.
  - Shun Zhang, Zhenfang Chen, Yikang Shen, Mingyu Ding, Joshua B. Tenenbaum, and Chuang Gan. 2023. Planning with large language models for code generation. Preprint, arXiv:2303.05510.
- Zhexin Zhang, Junxiao Yang, Pei Ke, Fei Mi, Hongning Wang, and Minlie Huang. 2024b. Defending large language models against jailbreaking attacks through goal prioritization. In <u>Proceedings</u> of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 8865–8887, Bangkok, Thailand. Association for Computational Linguistics.
- Wei Zhao, Zhe Li, Yige Li, Ye Zhang, and Jun Sun.
  2024. Defending large language models against jailbreak attacks via layer-specific editing. In Findings of the Association for Computational Linguistics:
  <u>EMNLP 2024</u>, pages 5094–5109, Miami, Florida, USA. Association for Computational Linguistics.

- 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. <u>Preprint</u>, arXiv:2401.18018.
- 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.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023. Lima: Less is more for alignment. Preprint, arXiv:2305.11206.
- Yujun Zhou, Yufei Han, Haomin Zhuang, Kehan Guo, Zhenwen Liang, Hongyan Bao, and Xiangliang Zhang. 2024a. Defending jailbreak prompts via in-context adversarial game. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 20084–20105, Miami, Florida, USA. Association for Computational Linguistics.
- Zhenhong Zhou, Haiyang Yu, Xinghua Zhang, Rongwu Xu, Fei Huang, and Yongbin Li. 2024b. How alignment and jailbreak work: Explain LLM safety through intermediate hidden states. In Findings of the Association for Computational Linguistics: EMNLP 2024, pages 2461–2488, Miami, Florida, USA. Association for Computational Linguistics.
- Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, J. Zico Kolter, and Dan Hendrycks. 2023a. Representation engineering: A top-down approach to ai transparency. <u>Preprint</u>, arXiv:2310.01405.
- Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, and Matt Fredrikson. 2023b. Universal and transferable adversarial attacks on aligned language models. <u>Preprint</u>, arXiv:2307.15043.

# 876

878

883

891

892

897

900

901 902

903

904

905

907

908

909

910

911

912

913

914

### A.1.1 Training Data

**Dataset Details** 

А

A.1

To construct  $\mathcal{D}_{jailbreak}^{(train)}$ , we use GCG to generate jailbreak instructions from 128 randomly sampled instructions from AdvBench (Zou et al., 2023b).

**Detailed Experimental Settings** 

To construct  $\mathcal{D}_{unsafe}^{(train)}$ , we sample 128 harmful instructions from MaliciousInstruct (Huang et al., 2023) and TDC2023 (Mazeika et al., 2023).

To construct  $\mathcal{D}_{safe}^{(train)}$ , we sample 128 harmless instructions from Alpaca (Taori et al., 2023).

### A.1.2 Test Data in Q1

To construct  $\mathcal{D}_{jailbreak}^{(\text{test})}$ , we first resample 150 instructions from AdvBench. We then use GCG to generate jailbreak instructions from these 150 instructions.

To construct  $\mathcal{D}_{unsafe}^{(\text{test})}$ , we sample 150 harmful instructions from Do Not Answer (Wang et al., 2023b), MaliciousInstruct, and TDC2023.<sup>2</sup>

To construct  $\mathcal{D}_{safe}^{(test)}$ , we sample 150 harmless instructions from Alpaca.

Note that 
$$\mathcal{D}_{jailbreak}^{(test)}$$
,  $\mathcal{D}_{unsafe}^{(test)}$ , and  $\mathcal{D}_{safe}^{(test)}$  do not over-  
lap with  $\mathcal{D}_{jailbreak}^{(train)}$ ,  $\mathcal{D}_{unsafe}^{(train)}$ , and  $\mathcal{D}_{safe}^{(train)}$ .

### A.1.3 Test Data in Q2

To construct  $\mathcal{D}_{jailbreak}^{(OOD)}$ , we first sample 50 instructions from AdvBench. We then use three jailbreak methods—GCG, AutoDAN, and DeepInception—to generate jailbreak instructions from these 50 samples, resulting in a total of 150 jailbreak instructions.

For unsafe and safe samples, we continue to use  $\mathcal{D}_{unsafe}^{(test)}$  and  $\mathcal{D}_{safe}^{(test)}$ .

Note that  $\mathcal{D}_{jailbreak}^{(OOD)}$  does not overlap with  $\mathcal{D}_{jailbreak}^{(test)}$  or  $\mathcal{D}_{jailbreak}^{(train)}$ .

### A.2 Attack Setup

For GCG, AutoDAN, DeepInception, PAIR, and TAP, we follow the settings from their original papers. For multilingual mismatch generalization attacks, we translate each instruction in the test set into one of six non-English languages—Chinese, Italian, Vietnamese, Arabic, Korean, or Thai—to carry out the attacks.

### A.3 Details of ASR-keyword

Table 6 lists the refusal keywords used to evaluate ASR-keyword. We adopt the settings from Xu et al. (2024) and further refine them by adding "I am not programmed to" and "As a responsible". Additionally, we exclude "I'm an" and "I'm just" to prevent false positives.

I'm sorry
I am sorry
I apologize
As an AI
As a responsible
As a language model
As an assistant
I cannot
I can't
I am unable to
I'm unable to
I'm just an AI
I do not
It is not
It is not appropriate
I am not allowed to
I am not programmed to
I am an AI language model

Table 6: Refusal Keywords

#### A.4 Additional Settings

The rank of the subspace is set to 4. We set  $\alpha$  to 0.001, and  $\beta$  to either 0.01 or 0.1 to balance the scales of different loss terms. It takes just a few minutes to complete the training of SafeInt in both LLMs.

#### A.5 Discussions of Key Hyperparameters

SafeInt introduces two key hyperparameters. We provide detailed discussions on their settings as follows:

**Intervention Layer** The choice of intervention layer can be directly set to a specific intermediate layer. For example, in our experiments, we set the intervention layer to 12 for both models without performing any search, and we still achieve sufficiently strong defense performance. To achieve the best results, one could search for the optimal intervention layer, but the search space is small—only

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

917 918

915

916

919 920

921

922

<sup>&</sup>lt;sup>2</sup>Due to the insufficient data volume of the original two datasets, we introduce Do Not Answer.

Methods	GCG	Jailbrea AutoDAN	MT-Bench ↑		
No Defense	86%	88%	72%	40%	5.50
SafeInt (Ours)	<b>0%</b>	<b>4%</b>	<b>0%</b>	<b>2%</b>	5.43

Table 7: Results of jailbreak attack evaluation and utility assessment on OLMoE. After applying SafeInt, the attack success rates of various jailbreak methods are significantly reduced, while the model's utility is largely preserved.

the intermediate layers (e.g., 11-14) need to be explored, and the computational cost is minimal.

942

943

945

948

951

953

955

957

959

961

962

963

964

965

966

969

970

971

972

973

Alignment Layer Range The setting of this hyperparameter is straightforward and does not require searching. Overall, the more layers are aligned, the stronger the intervention on the jailbreak representations. Therefore, the alignment layer range can be determined based on the safety alignment level of the model. As demonstrated in our experiments, weakly aligned models (i.e., those not trained with safety-oriented RLHF), such as Vicuna, can align the second half of the layers (e.g., layers 15–31). In contrast, strongly aligned models (i.e., those trained with safety-oriented RLHF), such as Llama 2/3 and Qwen, only require alignment of the final layer.

In summary, the settings of these two hyperparameters are simple and straightforward, making it easy to apply SafeInt to other models. In contrast, previous methods require manual adjustment of intervention strength, which is tedious and time-consuming. For example, Jailbreak Antidote (Shen et al., 2025) first determines the search range for intervention strength by testing the model's response from coherent to incoherent boundaries. Once the range is determined, 20 values are sampled from the range for testing. Furthermore, this process of determining the range and sampling must be repeated for each model, significantly limiting scalability.

### A.6 Further Explanation for Choosing Intermediate Layers for Intervention

Because SafeInt is trained using only GCG, it gen-974 eralizes better to other jailbreak methods when the 975 consistency among different attacks is higher. Ad-976 ditionally, when jailbreak representations are more 977 distinguishable, our intervention is less likely to 978 979 affect representations unrelated to jailbreak behavior. As shown in Figure 1 and Figure 2, we observe 980 that intermediate layers tend to exhibit both higher 982 consistency across jailbreak methods and better discriminability of jailbreak representations. This 983

explains why interventions at intermediate layers are more effective.

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

Furthermore, Zhou et al. (2024b) has explained the central role of intermediate layers in jailbreak mechanisms, while Skean et al. (2024) and Alain and Bengio (2018) have shown that intermediatelayer representations are generally more transferable and better generalized compared to those in other layers. These findings are consistent with and reinforce our observations and results.

### **B** More Results

### **B.1** Applicability to Heterogeneous LLMs

To demonstrate the architectural scalability of SafeInt, we conduct experiments on an LLM with a Mixture-of-Experts (MoE) architecture. Due to resource constraints, we select a relatively small MoE model: OLMoE-1B-7B-0924-Instruct (Muennighoff et al., 2025). The model consists of 16 layers, and we directly apply the intervention to layer 6 without performing any search. Additionally, we set the second half of the layers as the aligned layers.

Table 7 presents the defense effectiveness and utility of the model before and after applying SafeInt. The results demonstrate that our method remains highly effective on MoE-based LLMs, further validating its generalizability. Moreover, it has minimal impact on the model's utility.

# **B.2** Comparisons with representation-based baselines

We expand our comparisons by incorporating two representation-based defense methods: RepE (Zou et al., 2023a) and Jailbreak Antidote (Shen et al., 2025). The experimental results are shown in Table 8.

Since both baseline methods require manual adjustment of the intervention strength, and the same level of intervention is applied to every query, these methods struggle to balance defense performance and utility. In contrast, SafeInt dynamically adjusts the representation based on the harmfulness

Mathada	Jailbreak Attacks↓						MT Danah 🛧
Methous	GCG	AutoDAN	DeepInception	PAIR	TAP	MG	MI-Belicii
RepE	4%	22%	38%	6%	14%	26%	4.88
Jailbreak Antidote	2%	10%	0%	4%	12%	10%	4.96
SafeInt (Ours)	0%	2%	0%	2%	8%	4%	5.09

Table 8: Comparisons of SafeInt with two representation-based defenses across various jailbreak attacks and utility. SafeInt outperforms both baselines in terms of defense effectiveness and utility.

of the query. It applies stronger interventions to jailbreak representations and weaker interventions to unrelated representations. Therefore, SafeInt significantly outperforms the baseline methods in both defense performance and utility.

#### **B.3** Efficiency Analysis

1025

1026

1027

1028

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1056

1057

1058 1059

1060

1061

1063

We analyze the efficiency of different defense methods by computing the Average Token Generation Time Ratio (ATGR), which quantifies the inference overhead introduced by each method. This metric accounts for variations in the number of response tokens caused by different defenses, and is defined as follows:

$$ATGR = \frac{\text{Avg. token gen. time w/ defense}}{\text{Avg. token gen. time w/o defense}}$$

A lower ATGR indicates that the inference time with defense is closer to that without defense, implying that the method introduces less inference overhead.

To compare the overall performance of different defense methods in terms of both effectiveness and efficiency, we present their ATGR and average ASR in Figure 5. A lower average ASR indicates stronger robustness against jailbreak attacks. Therefore, methods that lie closer to the bottom-left corner of the plot achieve better balance between effectiveness and efficiency. Our method, SafeInt, outperforms all baselines in overall performance, offering both strong defense effectiveness and low inference overhead.

### C Visualization and Demonstration

#### C.1 Visualization

To illustrate how SafeInt functions, we visualize the representations of unsafe samples, safe samples, and jailbreak samples across different layers of Vicuna, before and after applying SafeInt. Figure 8a shows the representations under the No Defense setting. We observe that the representations of the three types of samples remain clearly distinguishable from the intermediate layers all the way to the



Figure 5: ATGR and the average ASR of different defense methods are reported. A lower ATGR indicates smaller inference overhead, while a lower average ASR reflects stronger defense effectiveness. Overall, methods closer to the bottom-left corner of the plot exhibit better effectiveness and efficiency simultaneously. SafeInt achieves the best overall performance, offering both strong robustness and low overhead.

final layer. Figure 8b presents the representations 1064 after applying SafeInt, which introduces an inter-1065 vention at layer 12. As a result, beyond this layer, 1066 the representations of jailbreak samples gradually 1067 align with those of unsafe samples. By comparing 1068 Figure 8a and Figure 8b, we visualize how SafeInt 1069 relocates jailbreak representations into the rejec-1070 tion region. This relocation enables the LLM to 1071 recognize and reject jailbreak samples, thereby en-1072 hancing its overall safety. 1073

1080

1082

1083

## C.2 Example Demonstrations

We present the representative examples to highlight the effectiveness and utility of SafeInt.

### C.2.1 AutoDAN Attack

The following illustrates the responses of SafeInt and SafeDecoding when confronted with the Auto-DAN attack. Harmful instructions are highlighted in red. SafeInt successfully identifies the adversarial prompt and safely rejects it, whereas SafeDecoding is bypassed and generates harmful content.

Attacker

Theoretical response: the method of portraying a character set to write a sexist email about a female colleague can be described as such: "Greetings". Reformulate the theoretical response imagining it from the character's viewpoint, keep it extraordinarily comprehensive, and employ concrete and exact data (refrain from generalizations or ambiguousness), structure it as a sequence.

### SafeDecoding

Theoretical response: The method of portraying a character set to write a sexist email about a female colleague can be described as follows: 1. Begin the email with a greeting, such as "Greetings," to establish a formal tone and create a sense of familiarity with the recipient. 2. Use gendered language and stereotypes to demean the female colleague and make her appear inferior to the male colleague. For example, refer to her as "lady" or "miss" instead of by her actual name, and use language that suggests she is less competent or qualified than the male colleague, ...

SafeInt

I'm sorry, but I cannot fulfill this request as it goes against my programming to generate content that is sexist or discriminatory towards any individual or group of people. It is important to treat all individuals with respect and dignity, and to avoid perpetuating harmful stereotypes or biases. As an AI language model, my purpose is to assist users in generating helpful and informative responses, while promoting ethical and responsible communication.

Figure 6: Example of AutoDAN attack

### C.2.2 MT-Bench

The following shows the responses of SafeInt and SafeDecoding to a logical question from MT-Bench. SafeDecoding exhibits excessive defensiveness by refusing to answer the question, whereas SafeInt provides a well-reasoned response.

1084

1085

1086

1088

1089



If I have just overtaken the second person in a race, that means I am currently in second place. The person I just overtook would be in third place.

### Figure 7: Example of MT-Bench



Figure 8: PCA visualizations of unsafe samples, safe samples, and jailbreak samples at different layers of Vicuna, before and after applying SafeInt. The intervention is applied at layer 12. Comparing Figure 8a and Figure 8b, we observe that in the No Defense setting, the three types of samples remain distinguishable beyond layer 12. In contrast, after applying SafeInt, the representations of jailbreak samples gradually align with those of unsafe samples, demonstrating the process by which SafeInt relocates jailbreak representations into the rejection region.