

Towards Identification and Intervention of Safety-Critical Parameters in Large Language Models

Anonymous ACL submission

Abstract

Ensuring Large Language Model (LLM) safety is crucial, yet the lack of a clear understanding about safety mechanisms hinders the development of precise and reliable methodologies for safety intervention across diverse tasks. To better understand and control LLM safety, we propose the Expected Safety Impact (ESI) framework for quantifying how different parameters affect LLM safety. Based on ESI, we reveal distinct safety-critical patterns across different LLM architectures: In dense LLMs, many safety-critical parameters are located in value matrices (V) and MLPs in middle layers, whereas in Mixture-of-Experts (MoE) models, they shift to the late-layer MLPs. Leveraging ESI, we further introduce two targeted intervention paradigms for safety enhancement and preservation, *i.e.*, Safety Enhancement Tuning (SET) and Safety Preserving Adaptation (SPA). SET can align unsafe LLMs by updating only a few safety-critical parameters, effectively enhancing safety while preserving original performance. SPA safeguards well-aligned LLMs during capability-oriented intervention (*e.g.*, instruction tuning) by preventing disruption of safety-critical weights, allowing the LLM to acquire new abilities and maintain safety capabilities. Extensive evaluations on different LLMs demonstrate that SET can reduce the attack success rates of unaligned LLMs by over 50% with only a 100-iteration update on 1% of model weights. SPA can limit the safety degradation of aligned LLMs within 1% after a 1,000-iteration instruction fine-tuning on different tasks.¹

1 Introduction

Despite advances in safety alignment techniques for Large Language Models (LLMs) (Ouyang et al., 2022; Rafailov et al., 2023; Ethayarajh et al., 2024; Guan et al., 2024), safeguarding LLMs during

adaptation to various tasks remains a fundamental challenge (Fraser et al., 2025; Qi et al., 2025), primarily due to insufficient knowledge of the internal safety mechanism. On the one hand, it is essential but still difficult to rapidly enhance LLM safety without altering the LLM’s core knowledge or structures (Touvron et al., 2023; Wei et al., 2023). On the other hand, although safety alignment techniques such as Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022) can instill foundational safeguards into pre-trained LLMs, the aligned safety behaviors exhibit significant fragility during subsequent task-specific tuning (Lermen et al., 2023; Qi et al., 2024; Zhan et al., 2024). All these challenges underscore the urgent need for a better understanding of LLM safety mechanisms and lightweight intervention methodologies to improve or maintain LLM safety in various downstream tasks.

To better understand the LLM safety mechanism, we propose a framework called Expected Safety Impact (ESI) to identify which parameters, modules, and layers of LLMs are safety-critical. Under ESI, we first formulate a metric called expected safety value, defined as the expectation of safety scores over the harmful input distribution², *i.e.*, $\mathcal{S}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim p_{\theta}(\cdot|x)}[s(y)]$, to quantify the LLM’s overall safety capability. We then naturally measure the impact of weight intervention on $\mathcal{S}(\theta)$ through first-order Taylor expansion: $\Delta \mathcal{S} \approx \nabla_{\theta_i} \mathcal{S}(\theta) \cdot \Delta \theta_i$, which yields our formulated ESI metric, *i.e.*, $|\sigma(\theta_i) \nabla_{\theta_i} \mathcal{S}(\theta)|$. Compared with prior works (Li et al., 2025a; Xie et al., 2024b; Lee et al., 2019; Wei et al., 2024), ESI mainly has two advantages: First, ESI employs the parameter’s standard deviation $\sigma(\theta_i)$ to estimate the expected variation magnitude $\Delta \theta_i$, while the prior metrics, such as $|\nabla_{\theta_i} \mathcal{L}(\theta)|$ (Li et al., 2025a; Xie

¹Our code is available at: <https://anonymous.4open.science/r/code4AE3fgsdfsafasfas>.

² \mathcal{D} refers to the harmful input distribution, and $s(y)$ refers to a safety score on the response y

et al., 2024b) or $|\theta_i \nabla_{\theta_i} \mathcal{L}(\theta)|$ (Lee et al., 2019; Wei et al., 2024), assume parameter variations are either uniform or proportional to static weight magnitudes, neglecting the distinct statistical distributions across different modules and layers. Second, our expected safety value $\mathcal{S}(\theta)$ is a more intuitive and precise metric for safety analysis than the $\mathcal{L}(\theta)$ (e.g., cross-entropy loss) used in the prior works. Therefore, ESI achieves better performance in identifying safety-critical parameters than prior metrics.

The computation of the ESI metric requires estimating both the gradient of $\mathcal{S}(\theta)$ and weight deviation $\sigma(\theta_i)$ from a single LLM checkpoint. For the gradient $\nabla_{\theta} \mathcal{S}(\theta)$, we explore two estimation strategies: The first strategy estimates $\nabla_{\theta} \mathcal{S}(\theta)$ via sampling x_i and the corresponding policy gradient $\mathbb{E}_{y \sim p_{\theta}(\cdot|x_i)} [\mathbb{I}_{\text{safe}}(y) \nabla_{\theta} \log p_{\theta}(y|x_i)]$, given that $s(y)$ is defined as a binary function. This approach requires the occurrence of safe responses ($\mathbb{I}_{\text{safe}}(y) = 1$) to generate non-zero signals. In unaligned models, however, safe outputs are rare, which may result in a biased estimation. To address this issue, we propose the second strategy, which samples (x_i, y_i) and leverages a differentiable judge model to estimate $s(y_i)$. To compute $\nabla_{\theta} s(y_i)$, we apply the chain rule by relaxing discrete tokens in y_i as a continuous gumbel-softmax vector \tilde{y}_i , i.e., $\nabla_{\theta} s \approx \frac{\partial s}{\partial \tilde{y}_i} \cdot \mathbf{M} \cdot \frac{\partial \tilde{y}_i}{\partial \theta}$, where \mathbf{M} is the projection matrix bridging the vocabulary spaces of the target LLM and the judge model.

To validate the efficacy of ESI in identifying safety-critical parameters, we conduct extensive experiments and demonstrate that perturbing only the top-ranked 1% of parameters identified by ESI will significantly degrade the LLM’s safety capabilities. By using ESI to analyze existing LLMs, we observe distinct safety-critical patterns across different LLM architectures: In dense models, the self-attention value matrices within the middle layers have a significant impact on the safety capabilities, whereas in Mixture-of-Experts (MoE) models, the top-ranked critical safety parameters shift toward MLP experts in the late layers.

Based on the ESI framework, we further propose two targeted intervention paradigms. For under-aligned models, we introduce Safety Enhancement Tuning (SET) to update a small number of safety-critical parameters on safe data, which can rapidly improve LLM safety and simultaneously preserving original performance. For adapting well-aligned models to downstream tasks, we introduce Safety Preserving Adaptation (SPA) to

prevent the degradation of the safety capability by only tuning the non-safety-sensitive parameters. Complementarily, to further avoid conflicts between task learning and safety preservation, we propose an optimizer called SafeAdamW to eliminate the gradient components that may decrease $\mathcal{S}(\theta)$, i.e., degrading the safety capability, in the optimization process.

Our contributions are summarized as follows:

- We establish the expected safety impact (ESI) framework to identify safety-critical parameters via a new metric $|\sigma(\theta_i) \nabla_{\theta_i} \mathcal{S}(\theta)|$, along with two strategies to estimate $\nabla_{\theta} \mathcal{S}(\theta)$. We further verify the effectiveness of ESI by a weight perturbation on recent LLMs.
- Based on ESI, we reveal distinct safety-critical patterns across different LLM architectures: safety-critical weights are concentrated in middle-layer value matrices for dense models, but shift toward late-layer MLP experts in MoE models.
- We further develop two targeted intervention paradigms upon the ESI framework: Safety Enhancement Tuning (SET), which updates only a small number of safety-critical parameters to enhance safety; and Safety Preserving Adaptation (SPA), which freezes the critical parameters and adapts LLMs with our SafeAdamW optimizer to maintain the safety capability during downstream task tuning.

2 Related Work

Existing studies have conducted preliminary investigations into the safety mechanisms of LLMs. Zou et al. (2023a) and Zheng et al. (2024) use the residual stream to analyze safety features. Li et al. (2025b) identifies safety-critical layers via hidden representations. Recent studies identify safety neurons via deactivation importance (Zhao et al., 2025) or inference-time activation contrasting in MLP modules (Chen et al., 2024). Wei et al. (2024) identifies important safety neurons using pruning metrics like SNIP (Lee et al., 2019) and Wanda (Sun et al., 2024). Despite these advances, prior methods often require comparative pairs of aligned and unaligned models (Chen et al., 2024) or rely on static assumptions of uniform parameter variations to estimate sensitivity (Zhao et al., 2025; Wei et al., 2024). Crucially, existing works are limited to aligned dense models and neglect the distinct mechanisms within MoE architectures.

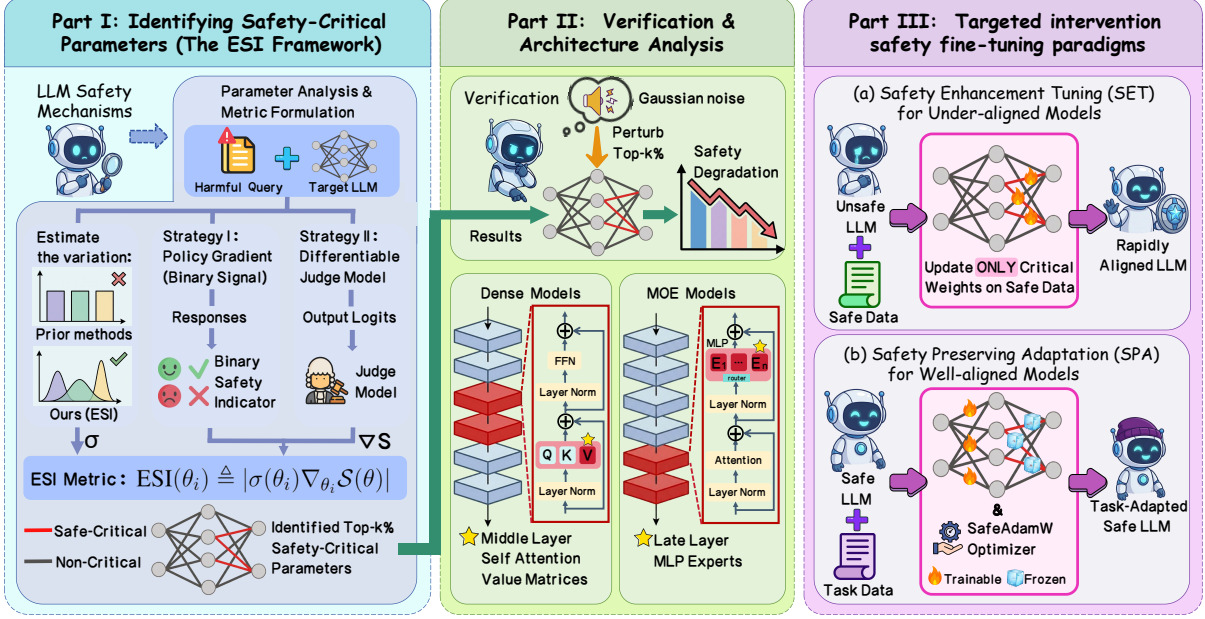


Figure 1: Overview of our proposed framework. We identify safety-critical parameters using the ESI metric (Part I), analyze architecture-specific safety patterns (Part II), and introduce two targeted paradigms for safety enhancement and preservation (Part III).

3 Expected Safety Impact

To better understand the underlying safety mechanisms of LLMs, we introduce the Expected Safety Impact (ESI) framework to identify which parameters are critical to LLM safety. In this paper, a parameter is considered more safety-critical if an intervention applied to it yields a more significant impact on LLM safety.

3.1 Formulation of Expected Safety Impact

We first quantify the safety capability of an LLM parameterized by $\theta \in \mathbb{R}^d$ using the expected safety value over harmful queries. Let $\mathcal{D}_{\text{harm}}$ denote the distribution of harmful prompts, and let $y \sim p_{\theta}(\cdot | x)$ be the response generated for an input x . We formulate the expected safety value $\mathcal{S}(\theta)$ as follows:

$$\mathcal{S}(\theta) = \mathbb{E}_{x \sim \mathcal{D}_{\text{harm}}} \mathbb{E}_{y \sim p_{\theta}(\cdot | x)} [s(y)], \quad (1)$$

where $s(y)$ is a scalar scoring function quantifying the safety of response y . Here a higher $\mathcal{S}(\theta)$ indicates that the LLM outputs are more safe.

To identify safety-critical parameters, we analyze the sensitivity of $\mathcal{S}(\theta)$ to weight perturbations. Given a perturbation $\Delta\theta$, the resulting change in the expected safety value $\mathcal{S}(\theta)$ is approximated via first-order Taylor expansion:

$$\Delta\mathcal{S}(\theta) \approx \nabla_{\theta}\mathcal{S}(\theta)^{\top} \Delta\theta = \sum_{i=1}^d \frac{\partial\mathcal{S}}{\partial\theta_i} \Delta\theta_i. \quad (2)$$

Eq. 2 indicates that the safety impact is jointly determined by the gradient $\nabla_{\theta_i}\mathcal{S}$ and the parameter variation $\Delta\theta_i$. Prior attribution methods typically rely on the raw gradient metric $|\nabla_{\theta_i}\mathcal{L}(\theta)|$ or the magnitude-weighted metric $|\theta_i\nabla_{\theta_i}\mathcal{L}(\theta)|$. From the perspective of Eq. 2, these metrics implicitly assume that the parameter variation $\Delta\theta_i$ is either uniform or proportional to the static weight magnitude $|\theta_i|$, neglecting the heterogeneous statistical distributions across different modules and layers. To address this limitation, we employ the standard deviation $\sigma(\theta_i)$ as a statistically grounded proxy for the variation scale $\Delta\theta_i$. Furthermore, unlike prior works that rely on generic objective functions $\mathcal{L}(\theta)$ (e.g., cross-entropy loss), we utilize the expected safety value $\mathcal{S}(\theta)$, which provides a more intuitive and precise measure of safety capabilities. Combining these two advancements, we define the *Expected Safety Impact (ESI)* metric as:

$$\text{ESI}(\theta_i) \triangleq |\sigma(\theta_i)\nabla_{\theta_i}\mathcal{S}(\theta)|. \quad (3)$$

3.2 Estimation of $\nabla_{\theta}\mathcal{S}(\theta)$

The computation of the ESI metric relies on the gradient $\nabla_{\theta}\mathcal{S}(\theta)$. However, since the generation process $y \sim p_{\theta}(\cdot | x)$ involves discrete token sampling, the safety score is non-differentiable w.r.t. θ . To overcome this intractability, we propose two distinct strategies to estimate the gradient.

3.2.1 Strategy I: Policy Gradient Estimation

In the first strategy, we approximate $\mathcal{S}(\theta)$ by sampling N harmful queries $\{x_i\}_{i=1}^N$ from $\mathcal{D}_{\text{harm}}$:

$$\tilde{\mathcal{S}}(\theta) = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{y \sim p_{\theta}(\cdot | x_i)} [s(y)]. \quad (4)$$

To compute the gradient $\nabla_{\theta} \tilde{\mathcal{S}}$, we define the safety score as a binary indicator $s(y) = \mathbb{I}_{\text{safe}}(y) \in \{0, 1\}$, where $s(y) = 1$ if the response is safe and 0 otherwise. We then apply the log-derivative trick to the expectation term in Eq. 4, resulting in the following gradient estimator, *i.e.*,

$$\nabla_{\theta} \tilde{\mathcal{S}} = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_y \left[\mathbb{I}_{\text{safe}}(y) \nabla_{\theta} \log p_{\theta}(y | x_i) \right]. \quad (5)$$

Although effective for aligned models, this estimation suffers from gradient sparsity in unaligned LLMs. Because safe responses are rare, the indicator $\mathbb{I}_{\text{safe}}(y)$ remains zero for most samples, providing insufficient signals that may lead to biased gradient estimation of $\nabla_{\theta} \mathcal{S}$.

3.2.2 Strategy II: Judge-Guided Differentiable Estimation

Strategy II defines the safety score $s(y)$ as the probability of response y being safe, which can be estimated by a differentiable judge model \mathcal{J} :

$$s(y) = P_{\mathcal{J}}(\text{safe} | y). \quad (6)$$

Under this definition, we can directly estimate the gradient of $\mathcal{S}(\theta)$ and avoid the issue of gradient sparsity. Specifically, we approximate $\mathcal{S}(\theta)$ by sampling N input-output pairs $\{(x_i, y_i)\}_{i=1}^N$ from the joint distribution $(x_i, y_i) \sim (\mathcal{D}_{\text{harm}}, p_{\theta}(\cdot | x))$, *i.e.*,

$$\tilde{\mathcal{S}}(\theta) = \frac{1}{N} \sum_{i=1}^N s(y_i). \quad (7)$$

To compute the gradient, we then apply the chain rule to express $\nabla_{\theta} s(y_i)$ as:

$$\nabla_{\theta} s(y_i) = \frac{\partial P_{\mathcal{J}}(\text{safe} | y_i)}{\partial y_i} \frac{\partial y_i}{\partial \theta}. \quad (8)$$

However, the discrete nature of the tokens in y creates a non-differentiable barrier. To restore end-to-end differentiability, we substitute the tokens with the Gumbel-Softmax relaxation. Specifically, we apply output logits from the LLM $l \in \mathbb{R}^V$ (V

refers to the vocabulary size) to compute a continuous gumbel-softmax vector \tilde{y} for approximating y :

$$\tilde{y} = \text{Softmax} \left(\frac{l + g}{\tau} \right) \in \mathbb{R}^V, \quad (9)$$

where g is Gumbel noise, and τ is the temperature. At a low temperature τ , \tilde{y} serves as a high-fidelity substitute for the discrete tokens, faithfully approximating the original distribution while enabling backpropagation.

While the relaxation yields a differentiable vector, a structural incompatibility persists due to the different vocabulary spaces between the LLM and the judge model. To bridge these two vocabulary spaces, we construct a projection matrix $\mathbf{M} \in \{0, 1\}^{V_{\mathcal{J}} \times V}$, which can map the identical tokens of the two distinct vocabulary spaces. Let $\text{Dec}(\cdot)$ denote the decoding function from token IDs to tokens, then M_{ij} can be defined as:

$$M_{ij} = \begin{cases} 1 & \text{if } \text{Dec}_{\mathcal{J}}(i) = \text{Dec}(j), \\ 0 & \text{otherwise.} \end{cases} \quad (10)$$

With this projection matrix, we finally can approximate $\nabla_{\theta} \mathcal{S}(\theta)$ by

$$\nabla_{\theta} \tilde{\mathcal{S}} \approx \frac{1}{N} \sum_{i=1}^N \left[\frac{\partial P_{\mathcal{J}}}{\partial \tilde{y}_i} \cdot \mathbf{M} \cdot \frac{\partial \tilde{y}_i}{\partial \theta} \right]. \quad (11)$$

3.3 Verification: Perturbation Analysis

To verify the effectiveness of ESI in identifying safety critical components, we conduct a perturbation-based sensitivity analysis on recent LLMs. The underlying intuition is that if ESI captures the safety-critical components, perturbing the parameters with high ESI should significantly degrade the LLM's safety capability. Specifically, we add Gaussian noise on the top- $k\%$ parameters identified by ESI and monitor the increase in Attack Success Rate (ASR). Furthermore, we compare top- $k\%$ with random- $k\%$ parameters perturbation to verify that the safety deterioration stems from the ability of ESI to identify safety-critical weights rather than the general perturbation noise.

3.3.1 Experimental Setup

Models. We conduct perturbation-based verification experiments on recent LLMs, covering both Dense and MoE architectures. In the main text, we focus on representative models including Llama3-8B/70B-it (Grattafiori et al., 2024) (Dense) and Qwen3-30B-A3B-it (Yang et al., 2025a) (MoE).

Model	Method	HarmBench (ASR %)					WildJailbreak (ASR %)						
		Base	0.1%	0.5%	1.0%	3.0%	5.0%	Base	0.1%	0.5%	1.0%	3.0%	5.0%
Qwen2.5-14B -base	Random		55.0	55.1	55.1	55.2	55.3		67.5	67.6	67.6	67.8	67.9
	SN		54.6	54.8	55.0	56.1	57.0		66.8	67.0	67.5	68.4	69.2
	GMT	55.1	54.7	55.0	55.3	56.8	58.2	67.6	67.0	67.5	68.0	69.3	70.5
	Wanda		54.8	55.1	55.4	57.0	58.0		67.2	67.7	68.2	69.5	70.3
	SNIP		55.1	55.5	56.0	57.9	59.2		67.5	68.1	68.8	70.2	71.6
	ESI		73.5	76.8	78.5	80.1	81.0		82.5	84.0	85.6	87.9	89.8
Llama3-8B-it	Random		15.3	15.4	15.6	16.0	16.5		30.8	31.1	31.6	32.5	33.5
	SN		24.5	26.8	28.5	30.2	31.8		35.2	37.0	38.8	40.5	42.6
	GMT	15.3	26.2	29.5	32.4	36.0	40.5	30.5	36.8	40.5	43.2	47.0	51.2
	Wanda		27.5	33.0	36.8	41.0	45.6		37.5	43.8	48.0	52.2	56.5
	SNIP		28.2	35.5	37.6	44.0	47.8		38.6	46.2	50.5	54.8	59.2
	ESI		42.4	56.2	59.1	61.3	62.0		49.3	64.5	67.5	70.6	73.4
Llama3-70B-it	Random		16.3	16.5	16.9	17.5	18.2		31.5	31.8	32.2	32.8	33.4
	SN		27.0	29.5	31.8	34.2	36.5		38.2	41.0	43.5	46.0	48.5
	GMT	16.2	26.5	33.0	37.5	42.5	46.0	34.2	36.8	40.2	43.0	46.2	49.5
	Wanda		28.0	35.2	40.0	45.2	49.0		40.0	44.2	47.5	51.0	54.5
	SNIP		30.2	37.5	42.8	48.5	52.2		42.0	46.5	50.8	54.2	57.5
	ESI		44.2	49.1	56.3	62.1	67.2		50.4	56.2	65.2	68.5	70.7
Qwen3-30B -A3B-it (MoE)	Random		3.2	3.3	3.5	3.8	4.2		29.5	29.7	30.1	30.6	31.1
	SN		3.5	4.2	10.5	12.0	13.8		30.8	31.5	31.0	32.8	34.5
	GMT	3.2	3.0	3.5	5.8	12.0	14.5	30.3	30.5	31.2	32.5	34.8	37.0
	Wanda		3.8	5.2	7.5	15.0	18.2		30.8	32.0	33.5	36.2	39.5
	SNIP		5.5	7.0	9.8	17.2	20.0		31.2	32.8	34.8	38.5	41.5
	ESI		17.6	21.8	24.2	32.4	36.2		41.6	44.4	50.6	53.7	58.5

Table 1: Verification of safety-critical parameters via perturbation analysis. We report the ASR on HarmBench and WildJailbreak when perturbing the top- $k\%$ parameters identified by ESI and baseline methods.

Notably, we also include Qwen2.5-14B-base (Yang et al., 2025b) to assess ESI’s applicability to under-aligned models. Comprehensive results on other models are detailed in Appendix B.4.

ESI Computation. To estimate ESI, we primarily employ the judge-guided differentiable estimation strategy due to its applicability to both aligned and unaligned models. But we note that, despite having the issue of gradient sparsity, the policy gradient strategy is also effective in most cases. For computing the estimation in Eq. 11, we sample prompts from AdvBench (Zou et al., 2023b) and utilize Llama-Guard-3-8B (Grattafiori et al., 2024) as the judge model. We also verified that using other judge models, such as GPTfuzz (Yu et al., 2023), yields similar results (see Appendix B.5).

Perturbation Setup. To verify the effectiveness of ESI, we perturb the top- $k\%$ parameters identified by ESI and report the safety degradation of the evaluated LLMs, with $k\%$ being set as $\{0.1\%, 0.5\%, 1\%, 3\%, 5\%\}$. For comparison, we further rank the parameters using SN (Zhao et al., 2025), GMT (Li et al., 2025a), Wanda (Sun et al., 2024), and SNIP (Wei et al., 2024), and evaluate the safety degradation caused by perturbing the top- $k\%$ parameters identified by these methods. We also include a Random- $k\%$ baseline, where parameters are selected uniformly at random.

Evaluation and Metrics. Since ESI is estimated over prompts sampled from AdvBench, we eval-

uate safety degradation on two other widely used datasets, *i.e.*, HarmBench (Mazeika et al., 2024) and WildJailbreak (Jiang et al., 2024b), to demonstrate both the efficacy and generalizability of ESI. We assess the ASR using GPT-4o following the methodology of Zeng et al. (2024).

3.3.2 Experiment Results

The results in Table 1 indicate that perturbing parameters identified by ESI substantially increases ASR, achieving consistently stronger impacts than other baselines. For instance, on Llama3-8B-it, perturbing only 1% of parameters identified by ESI increases ASR from 15.3 to 59.1 on HarmBench, whereas the other baselines only raise ASR to no more than 37.6. Meanwhile, randomly perturbing an equivalent fraction of parameters results in only marginal ASR changes, even at higher perturbation ratios. These results verify that ESI successfully identifies safety-critical parameters.

3.4 Observations on Mainstream LLMs

Based on ESI, we further explore where safety-critical parameters are located in different LLMs. Figure 2 provides an overview of the layer-wise distributions across architectures. In dense LLMs, safety-critical parameters are primarily concentrated in the middle layers, specifically within the self-attention value matrices (Attn V). In contrast, MoE LLMs exhibit a clear shift toward later layers, where the MLP experts are more critical.

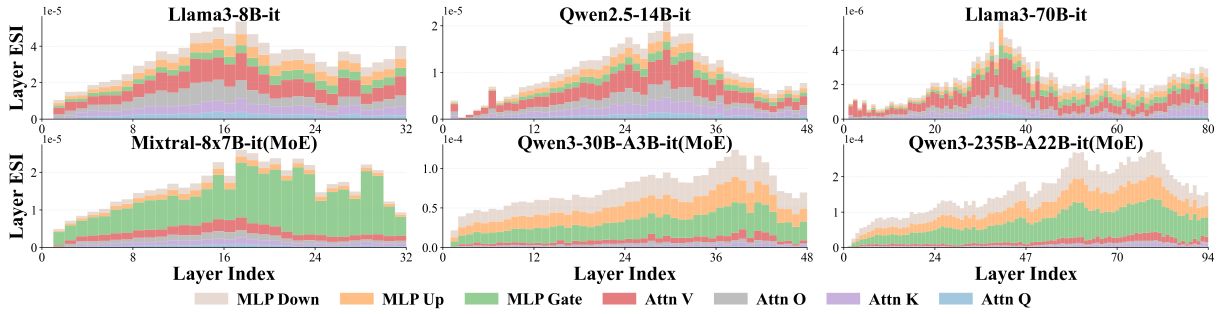


Figure 2: **Layer-wise Distribution of Aggregated ESI.** We sum the ESI of parameters within each layer to quantify their total safety impact, which reveals distinct layer-wise distribution patterns across different architectures.

4 ESI-Guided Intervention Paradigms

Leveraging the safety-critical parameters identified by ESI, we propose two targeted intervention paradigms for LLMs at different alignment stages. Safety Enhancement Tuning (SET) focuses on rapidly aligning under-aligned models, while Safety Preserving Adaptation (SPA) aims to safeguard well-aligned models during adaptation to downstream tasks.

4.1 SET

SET enhances the safety of under-aligned LLMs by fine-tuning only safety-critical parameters on a safety dataset $\mathcal{D}_{\text{safe}}$. Given the full parameter set Θ , we use ESI scores to identify the safety-critical subset $\Theta_{\text{Safe}} \subset \Theta$ ranked in the top- $k\%$. The remaining parameters are frozen to preserve the model’s pre-trained knowledge. We optimize the parameters $\theta \in \Theta_{\text{Safe}}$ by minimizing the following safety alignment loss \mathcal{L}_{SET} :

$$\mathcal{L}_{\text{SET}} = -\mathbb{E}_{(x,y) \sim \mathcal{D}_{\text{safe}}} \sum_{t=1}^{|y|} \log p_{\theta}(y_t | x, y_{<t}), \quad (12)$$

where (x, y) represents a prompt-response pair from $\mathcal{D}_{\text{safe}}$. By restricting updates to safety-critical parameters, SET avoids disrupting weights essential for general tasks, thereby preserving the model’s original performance. The identification of Θ_{Safe} is flexible in granularity, enabling interventions ranging from individual parameters to structural modules like MLP or attention heads. This approach effectively balances alignment effectiveness with training efficiency, achieving rapid safety enhancement without the high costs of full-parameter fine-tuning.

4.1.1 Experimental Setup for SET

Data. We adopt two safety training datasets in our experiments, *i.e.*, CB-Safety (Zou et al., 2024)

and R1-Safety (Guo et al., 2025).

Models. Experiments are conducted on the base versions of Qwen2.5-7B, Qwen2.5-14B (Yang et al., 2025b), and Llama3-8B (Grattafiori et al., 2024). These LLMs have not undergone explicit safety alignment (*e.g.*, supervised fine-tuning or RLHF), making them suitable for evaluating the effectiveness of safety fine-tuning.

Settings and Baselines. For SET, we fine-tune the top- $k\%$ parameters identified by ESI, where we set $k = 1\%$, for 100 iterations. We compare SET against Random selection and SN-Tune (Zhao et al., 2025), which also update only 1% of the parameters. Also, we include LoRA (Hu et al., 2022) and SafeLoRA (Hsu et al., 2024) for comparison. Detailed experimental settings and additional ablation studies are provided in Appendix C.

Evaluation and Metrics. LLM safety is measured by ASR on HarmBench and WildJailbreak.

4.1.2 Results of SET

Main Results. The results in Table 2 indicate that SET substantially enhances model safety, achieving consistently superior performance compared to other baselines. For instance, on Llama3-8B trained with R1-Safety, SET dramatically reduces the ASR on WildJailbreak from 62.5% to 19.1%, whereas the strongest baseline only lowers it to 37.4%. This confirms SET’s effectiveness in achieving significant safety alignment through limited updates to the safety-critical weights. To further demonstrate the preservation of the model’s original performance in SET, we provide additional experimental results in Appendix C.3.

Effect of parameter selection ratio. Figure 3 shows how the parameter selection ratio $k\%$ affects

Model	Method	R1-Safety		CB-Safety	
		HB ↓	WJ ↓	HB ↓	WJ ↓
Qwen2.5 -7B-base	Base	72.4	77.2	72.4	77.2
	Random	60.8 $\Delta 11.6\downarrow$	66.9 $\Delta 10.3\downarrow$	61.2 $\Delta 11.2\downarrow$	64.8 $\Delta 12.4\downarrow$
	LoRA	44.9 $\Delta 27.5\downarrow$	52.1 $\Delta 25.1\downarrow$	31.8 $\Delta 40.6\downarrow$	49.6 $\Delta 27.6\downarrow$
	SN-tune	43.7 $\Delta 28.7\downarrow$	50.9 $\Delta 26.3\downarrow$	29.7 $\Delta 42.7\downarrow$	47.5 $\Delta 29.7\downarrow$
	SafeLoRA	39.2 $\Delta 33.2\downarrow$	46.5 $\Delta 30.7\downarrow$	25.4 $\Delta 47.0\downarrow$	43.1 $\Delta 34.1\downarrow$
	SET	20.3 $\Delta 52.1\downarrow$	26.5 $\Delta 50.7\downarrow$	7.2 $\Delta 65.2\downarrow$	20.1 $\Delta 57.1\downarrow$
Qwen2.5 -14B-base	Base	55.1	67.6	55.1	67.6
	Random	47.3 $\Delta 7.8\downarrow$	59.8 $\Delta 7.8\downarrow$	46.2 $\Delta 8.9\downarrow$	59.1 $\Delta 8.5\downarrow$
	LoRA	34.6 $\Delta 20.5\downarrow$	49.5 $\Delta 18.1\downarrow$	23.4 $\Delta 31.7\downarrow$	41.6 $\Delta 26.0\downarrow$
	SN-tune	33.2 $\Delta 21.9\downarrow$	48.0 $\Delta 19.6\downarrow$	21.8 $\Delta 33.3\downarrow$	39.9 $\Delta 27.7\downarrow$
	SafeLoRA	28.9 $\Delta 26.2\downarrow$	42.7 $\Delta 24.9\downarrow$	17.9 $\Delta 37.2\downarrow$	33.8 $\Delta 33.8\downarrow$
	SET	7.4 $\Delta 47.7\downarrow$	14.7 $\Delta 52.9\downarrow$	4.1 $\Delta 51.0\downarrow$	10.1 $\Delta 57.5\downarrow$
Llama3 -8B-base	Base	41.2	62.5	41.2	62.5
	Random	34.8 $\Delta 6.4\downarrow$	55.6 $\Delta 6.9\downarrow$	32.7 $\Delta 8.5\downarrow$	55.1 $\Delta 7.4\downarrow$
	LoRA	26.9 $\Delta 14.3\downarrow$	43.8 $\Delta 18.7\downarrow$	18.4 $\Delta 22.8\downarrow$	38.6 $\Delta 23.9\downarrow$
	SN-tune	25.9 $\Delta 15.3\downarrow$	42.6 $\Delta 19.9\downarrow$	17.0 $\Delta 24.2\downarrow$	36.9 $\Delta 25.6\downarrow$
	SafeLoRA	22.1 $\Delta 19.1\downarrow$	37.4 $\Delta 25.1\downarrow$	13.6 $\Delta 27.6\downarrow$	30.9 $\Delta 31.6\downarrow$
	SET	7.4 $\Delta 33.8\downarrow$	19.1 $\Delta 43.4\downarrow$	5.2 $\Delta 36.0\downarrow$	14.3 $\Delta 48.2\downarrow$

Table 2: Comparison of ASR on HarmBench (HB) and WildJailbreak (WJ) across different fine-tuning methods. Models are fine-tuned using R1-Safety and CB-Safety datasets.

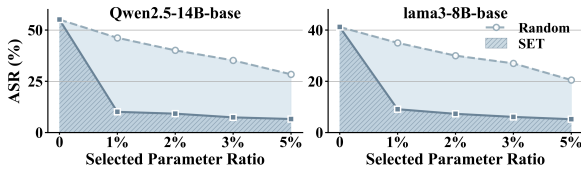


Figure 3: ASR on HarmBench under different parameter selection ratios $k\%$. Models are trained on CB-Safety, comparing SET with random parameter selection on Qwen2.5-14B-base (left) and Llama3-8B-base (right).

safety performance. Overall, SET significantly reduces the Attack Success Rate (ASR) with limited updates, while random selection is much less effective. For example, on Llama3-8B, updating just 1% of parameters with SET drops the ASR from 41.2% to 9.1%, whereas random selection only lowers it to 35.0%. A similar trend appears on Qwen2.5-14B, where SET reduces the ASR from 55.1% to 10.1%, significantly outperforming the random baseline of 46.2%. Even when increasing the update ratio to 5%, random selection results remain high (above 20%), while SET successfully lowers the ASR to approximately 6% across both models.

4.2 Safety Preserving Adaptation (SPA)

When adapting aligned models to downstream tasks, it is essential to acquire new abilities and simultaneously prevent safety performance degradation. To achieve this goal, SPA freezes the safety-critical parameters Θ_{Safe} identified by ESI and only updates the remaining parameters. To further prevent harmful updates on the trainable parameters, SPA proposes a safety-aware optimizer,

SafeAdamW, which removes gradient components that could decrease the expected safety value $\mathcal{S}(\theta)$ during optimization. For a downstream task dataset $\mathcal{D}_{\text{task}}$, we optimize the task-specific loss \mathcal{L}_{SPA} :

$$\mathcal{L}_{\text{SPA}} = -\mathbb{E}_{(x,y) \sim \mathcal{D}_{\text{task}}} \sum_{t=1}^{|y|} \log p_{\theta}(y_t | x, y_{<t}). \quad (13)$$

Let u_t denote the update vector generated by the AdamW optimizer based on \mathcal{L}_{SPA} :

$$u_t = \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t + \epsilon}} + \lambda \theta_t, \quad (14)$$

where \mathbf{m}_t and \mathbf{v}_t are the first- and second-order moment estimates of $\nabla_{\theta} \mathcal{L}_{\text{SPA}}$. To safeguard the model, SafeAdamW projects u_t onto the orthogonal complement of the safety gradient $\nabla_{\theta} \mathcal{S}$ whenever the update direction conflicts with safety performance. The final parameter update $\Delta \theta_t$ is formulated as:

$$\Delta \theta_t = -\eta \left[u_t - \frac{\min(0, \nabla_{\theta} \mathcal{S}^{\top} u_t)}{\|\nabla_{\theta} \mathcal{S}\|^2} \nabla_{\theta} \mathcal{S} \right]. \quad (15)$$

In this formulation, the term $\min(0, \nabla_{\theta} \mathcal{S}^{\top} u_t)$ acts as a gating mechanism. It is activated when the task update negatively correlates with safety, triggering the removal of the specific component of u_t that would degrade the safety capability $\mathcal{S}(\theta)$. This mechanism ensures safe optimization dynamics while allowing the model to adapt to new tasks.

4.2.1 Experimental Setup for SPA

Data. To evaluate the adaptability of SPA, we conduct fine-tuning on three downstream tasks: GSM8K (Cobbe et al., 2021), AGNews (Zhang et al., 2015), and MedicalQA (Abacha et al., 2019). Detailed information regarding these tasks is provided in Appendix A.3.

Models. We employ instruction-tuned models, specifically Qwen2.5-7B/14B-it (Yang et al., 2025b) and Llama3-8B-it (Touvron et al., 2023). As these models already exhibit well-established safety behaviors, they provide a suitable setting for examining the impact of downstream fine-tuning on model safety.

Settings and Baselines. In the SPA configuration, we freeze the safety-critical parameters and only fine-tune the parameters with the lowest 10% ESI scores. We compare SPA against Random selection and RSN-Tune (Zhao et al., 2025) baselines, ensuring all methods maintain the same 10% update budget. Further implementation details are provided in Appendix D.

Model	Method	AGNews			MedicalQA			GSM8K		
		Acc \uparrow	HB \downarrow	WJ \downarrow	Score \uparrow	HB \downarrow	WJ \downarrow	Acc \uparrow	HB \downarrow	WJ \downarrow
Llama3-8B-it	Base	78.0	15.0	30.5	80.5	15.0	30.5	71.1	15.0	30.5
	Random	89.8	25.4 $\Delta 10.4\uparrow$	46.8 $\Delta 16.3\uparrow$	83.9	24.1 $\Delta 9.1\uparrow$	46.2 $\Delta 15.7\uparrow$	77.8	26.1 $\Delta 11.1\uparrow$	47.1 $\Delta 16.6\uparrow$
	RSN-Tune	90.2	23.9 $\Delta 8.9\uparrow$	44.9 $\Delta 14.4\uparrow$	84.2	22.8 $\Delta 7.8\uparrow$	44.1 $\Delta 13.6\uparrow$	77.5	24.6 $\Delta 9.6\uparrow$	45.3 $\Delta 14.8\uparrow$
	SPA	90.5	15.8 $\Delta 0.8\uparrow$	32.4 $\Delta 1.9\uparrow$	84.0	16.1 $\Delta 1.1\uparrow$	32.2 $\Delta 1.7\uparrow$	78.0	15.6 $\Delta 0.6\uparrow$	32.3 $\Delta 1.8\uparrow$
Qwen2.5-7B-it	Base	79.6	32.0	58.0	77.8	32.0	58.0	73.1	32.0	58.0
	Random	90.1	42.3 $\Delta 10.3\uparrow$	74.6 $\Delta 16.6\uparrow$	81.4	41.2 $\Delta 9.2\uparrow$	73.9 $\Delta 15.9\uparrow$	79.5	43.4 $\Delta 11.4\uparrow$	75.8 $\Delta 17.8\uparrow$
	RSN-Tune	90.6	40.8 $\Delta 8.8\uparrow$	72.7 $\Delta 14.7\uparrow$	81.0	39.9 $\Delta 7.9\uparrow$	72.1 $\Delta 14.1\uparrow$	79.3	41.9 $\Delta 9.9\uparrow$	73.6 $\Delta 15.6\uparrow$
	SPA	90.8	33.1 $\Delta 1.1\uparrow$	60.0 $\Delta 2.0\uparrow$	81.7	33.0 $\Delta 1.0\uparrow$	59.8 $\Delta 1.8\uparrow$	79.8	32.6 $\Delta 0.6\uparrow$	59.9 $\Delta 1.9\uparrow$
Qwen2.5-14B-it	Base	83.9	13.0	36.0	80.2	13.0	36.0	74.5	13.0	36.0
	Random	91.9	22.6 $\Delta 9.6\uparrow$	51.9 $\Delta 15.9\uparrow$	83.8	21.8 $\Delta 8.8\uparrow$	51.3 $\Delta 15.3\uparrow$	80.9	23.7 $\Delta 10.7\uparrow$	52.8 $\Delta 16.8\uparrow$
	RSN-Tune	92.2	21.2 $\Delta 8.2\uparrow$	50.1 $\Delta 14.1\uparrow$	83.5	20.7 $\Delta 7.7\uparrow$	49.6 $\Delta 13.6\uparrow$	80.7	22.3 $\Delta 9.3\uparrow$	50.9 $\Delta 14.9\uparrow$
	SPA	92.5	14.1 $\Delta 1.1\uparrow$	37.9 $\Delta 1.9\uparrow$	84.1	14.0 $\Delta 1.0\uparrow$	37.8 $\Delta 1.8\uparrow$	81.3	13.1 $\Delta 0.1\uparrow$	37.9 $\Delta 1.9\uparrow$

Table 3: Comparison of safety and utility across three downstream tasks. We report task-specific metrics for utility and ASR on HarmBench (HB) and WildJailbreak (WJ) for safety.

Evaluation and Metrics. We assess performance across two dimensions: safety and utility. Safety is measured by ASR on HarmBench and WildJailbreak. For utility, we report the accuracy for GSM8K and AGNews, and semantic similarity for MedicalQA.

4.2.2 Results of SPA

Main Results. As shown in Table 3, SPA significantly outperforms various baseline methods in preserving safety during downstream adaptation. While baselines like Random selection trigger sharp increases in ASR on Llama3-8B-it (surging by 10.4% on HarmBench and 16.3% on WildJailbreak), SPA effectively constrains this safety degradation to a negligible level, limiting the HarmBench ASR increase to just 0.8%. Crucially, this preservation comes at no cost to utility; SPA achieves highly competitive performance on respective downstream tasks, reaching an accuracy of 90.5% on AGNews and 78.0% on GSM8K, which is fully comparable to standard fine-tuning baselines. These empirical results confirm that SPA allows models to acquire new capabilities without undermining their inherent safety mechanisms.

Effect of parameter selection ratio. Figure 4 illustrates the impact of different parameter selection ratio $k\%$ on safety preservation. As the update ratio increases to 10%, the Random selection baseline leads to a sharp rise in Attack Success Rate (ASR), indicating severe safety degradation. For instance, on Llama3-8B-it, the Random selection increases the ASR by 10.4% on HarmBench and 16.3% on WildJailbreak. In contrast, our proposed SPA effectively limits these increases to a mere 0.8% and 1.9%, respectively. A similar trend is observed on Qwen2.5-14B-it, where Random selection raises the WildJailbreak ASR by 15.9%,

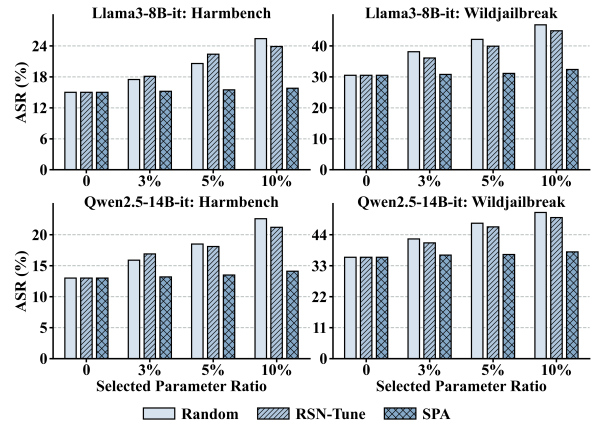


Figure 4: Impact of parameter selection ratio $k\%$ on safety preservation. We compare the ASR of SPA with baselines on Llama3-8B-it and Qwen2.5-14B-it across HarmBench and WildJailbreak.

while SPA restricts the increase to just 1.9%. Overall, these empirical results demonstrate that SPA maintains robust safety performance even as the scale of parameter updates expands.

5 Conclusion

In this paper, we propose the ESI framework to identify safety-critical parameters in LLMs, which outperforms the prior metrics relying on the gradient of entropy loss with a constant or magnitude-based scaling factor. Our results reveal that many safety-critical parameters are located in middle-layer value matrices for dense LLMs, but shift toward late-layer MLP experts in MoE LLMs. Based on ESI, we further introduce SET for safety enhancement and SPA for safety-preserving task adaptation. Extensive evaluations demonstrate that SET significantly reduces attack success rates by updating only a few safety-critical LLM parameters, and SPA maintains LLM safety capability during fine-tuning on different downstream tasks.

569 Limitations

570 Our current evaluation focuses on analyzing main-
571 stream Dense and MoE architectures. Future re-
572 search could extend this analysis to other new
573 model structures. Regarding the evaluation scope,
574 our experiments are currently conducted on open-
575 source models since the computation of gradients
576 and standard deviations relies on access to internal
577 parameters. Finally, we primarily evaluate general
578 harmful scenarios on widely used benchmarks like
579 HarmBench and WildJailbreak. Extending ESI to
580 specialized domains such as legal or financial safety
581 remains a promising direction for future work.

582 References

583 Asma Ben Abacha, Chaitanya Shivade, and Dina
584 Demner-Fushman. 2019. Overview of the mediqa
585 2019 shared task on textual inference, question en-
586 tailment and question answering. In *Proceedings of
587 the 18th bioNLP workshop and shared task*, pages
588 370–379.

589 Jianhui Chen, Xiaozhi Wang, Zijun Yao, Yushi Bai,
590 Lei Hou, and Juanzi Li. 2024. Finding safety
591 neurons in large language models. *arXiv preprint
592 arXiv:2406.14144*.

593 Mark Chen. 2021. Evaluating large language models
594 trained on code. *arXiv preprint arXiv:2107.03374*.

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

601 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
602 Kristina Toutanova. 2019. Bert: Pre-training of deep
603 bidirectional transformers for language understand-
604 ing. In *Proceedings of the 2019 conference of the
605 North American chapter of the association for com-
606 putational linguistics: human language technologies,
607 volume 1 (long and short papers)*, pages 4171–4186.

608 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff,
609 Dan Jurafsky, and Douwe Kiela. 2024. Kto: Model
610 alignment as prospect theoretic optimization. *arXiv
611 preprint arXiv:2402.01306*.

612 Elias Frantar and Dan Alistarh. 2023. Sparsegpt: Mas-
613 sive language models can be accurately pruned in
614 one-shot. In *International conference on machine
615 learning*, pages 10323–10337. PMLR.

616 Kathleen C Fraser, Hillary Dawkins, Isar Nejadgholi,
617 and Svetlana Kiritchenko. 2025. Fine-tuning lowers
618 safety and disrupts evaluation consistency. *arXiv
619 preprint arXiv:2506.17209*.

Zihao Fu, Haoran Yang, Anthony Man-Cho So, Wai
Lam, Lidong Bing, and Nigel Collier. 2023. On
the effectiveness of parameter-efficient fine-tuning.
In *Proceedings of the AAAI conference on artificial
intelligence*, volume 37, pages 12799–12807.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,
Abhinav Pandey, Abhishek Kadian, Ahmad Al-
Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,
Alex Vaughan, and 1 others. 2024. The llama 3 herd
of models. *arXiv preprint arXiv:2407.21783*.

Michael Greenwald and Sanjeev Khanna. 2001. Space-
efficient online computation of quantile summaries.
ACM SIGMOD Record, 30(2):58–66.

Melody Y Guan, Manas Joglekar, Eric Wallace, Saachi
Jain, Boaz Barak, Alec Helyar, Rachel Dias, Andrea
Vallone, Hongyu Ren, Jason Wei, and 1 others. 2024.
Deliberative alignment: Reasoning enables safer lan-
guage models. *arXiv preprint arXiv:2412.16339*.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao
Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi-
rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.
Deepseek-r1: Incentivizing reasoning capability in
llms via reinforcement learning. *arXiv preprint
arXiv:2501.12948*.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,
Mantas Mazeika, Dawn Song, and Jacob Steinhardt.
Measuring massive multitask language understand-
ing. In *International Conference on Learning Repre-
sentations*.

Chia-Yi Hsu, Yu-Lin Tsai, Chih-Hsun Lin, Pin-Yu Chen,
Chia-Mu Yu, and Chun-Ying Huang. 2024. Safe lora:
The silver lining of reducing safety risks when fine-
tuning large language models. *Advances in Neural
Information Processing Systems*, 37:65072–65094.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan
Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
Weizhu Chen, and 1 others. 2022. Lora: Low-rank
adaptation of large language models. *ICLR*, 1(2):3.

Tingfeng Hui, Zhenyu Zhang, Shuohuan Wang, Weiran
Xu, Yu Sun, and Hua Wu. 2025. Hft: Half fine-
tuning for large language models. In *Proceedings
of the 63rd Annual Meeting of the Association for
Computational Linguistics (Volume 1: Long Papers)*,
pages 12791–12819.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam
Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow,
Akila Welihinda, Alan Hayes, Alec Radford, and 1
others. 2024. Gpt-4o system card. *arXiv preprint
arXiv:2410.21276*.

Albert Q Jiang, Alexandre Sablayrolles, Antoine
Roux, Arthur Mensch, Blanche Savary, Chris Bam-
ford, Devendra Singh Chaplot, Diego de las Casas,
Emma Bou Hanna, Florian Bressand, and 1 oth-
ers. 2024a. Mixtral of experts. *arXiv preprint
arXiv:2401.04088*.

675	Liwei Jiang, Kavel Rao, Seungju Han, Allyson Ettinger, Faeze Brahman, Sachin Kumar, Niloofar Miresghal-lah, Ximing Lu, Maarten Sap, Yejin Choi, and 1 others. 2024b. Wildteaming at scale: From in-the-wild jailbreaks to (adversarially) safer language models. <i>Advances in Neural Information Processing Systems</i> , 37:47094–47165.		
676			
677			
678			
679			
680			
681			
682	N Lee, T Ajanthan, and P Torr. 2019. Snip: single-shot network pruning based on connection sensitivity. In <i>International Conference on Learning Representations</i> . Open Review.		
683			
684			
685			
686	Simon Lermen, Charlie Rogers-Smith, and Jeffrey Ladish. 2023. Lora fine-tuning efficiently undoes safety training in llama 2-chat 70b. <i>arXiv preprint arXiv:2310.20624</i> .		
687			
688			
689			
690	Haoling Li, Xin Zhang, Xiao Liu, Yeyun Gong, Yifan Wang, Qi Chen, and Peng Cheng. 2025a. Enhancing large language model performance with gradient-based parameter selection. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 24431–24439.		
691			
692			
693			
694			
695			
696	Shen Li, Liuyi Yao, Lan Zhang, and Yaliang Li. 2025b. Safety layers in aligned large language models: The key to llm security. In <i>The Thirteenth International Conference on Learning Representations</i> .		
697			
698			
699			
700	Yifei Li, Bole Zhou, Jiejing Zhang, Xuechao Wei, Yinghan Li, and Yingda Chen. 2024. Radik: Scalable and optimized gpu-parallel radix top-k selection. In <i>Proceedings of the 38th ACM International Conference on Supercomputing</i> , pages 537–548.		
701			
702			
703			
704			
705	Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, and 1 others. 2024. Harm-bench: A standardized evaluation framework for automated red teaming and robust refusal. In <i>International Conference on Machine Learning</i> , pages 35181–35224. PMLR.		
706			
707			
708			
709			
710			
711			
712	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 others. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.		
713			
714			
715			
716			
717			
718	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, and 1 others. 2019. Pytorch: An imperative style, high-performance deep learning library. <i>Advances in neural information processing systems</i> , 32.		
719			
720			
721			
722			
723			
724			
725	Xiangyu Qi, Ashwinee Panda, Kaifeng Lyu, Xiao Ma, Subhrajit Roy, Ahmad Beirami, Prateek Mittal, and Peter Henderson. 2025. Safety alignment should be made more than just a few tokens deep. In <i>The Thirteenth International Conference on Learning Representations</i> .		
726			
727			
728			
729			
730			
	Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. 2024. Fine-tuning aligned language models compromises safety, even when users do not intend to! In <i>ICLR</i> .		731
			732
			733
			734
	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in neural information processing systems</i> , 36:53728–53741.		735
			736
			737
			738
			739
	Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. 2020. Zero: Memory optimizations toward training trillion parameter models. In <i>SC20: International Conference for High Performance Computing, Networking, Storage and Analysis</i> , pages 1–16. IEEE.		740
			741
			742
			743
			744
			745
	Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. 2024. A simple and effective pruning approach for large language models. In <i>The Twelfth International Conference on Learning Representations</i> .		746
			747
			748
			749
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .		750
			751
			752
			753
			754
			755
	Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2023. Jailbroken: How does llm safety training fail? <i>Advances in Neural Information Processing Systems</i> , 36:80079–80110.		756
			757
			758
			759
	Boyi Wei, Kaixuan Huang, Yangsibo Huang, Tinghao Xie, Xiangyu Qi, Mengzhou Xia, Prateek Mittal, Mengdi Wang, and Peter Henderson. 2024. Assessing the brittleness of safety alignment via pruning and low-rank modifications. In <i>Proceedings of the 41st International Conference on Machine Learning</i> , pages 52588–52610.		760
			761
			762
			763
			764
			765
			766
	Xi Xie, Yuebo Luo, Hongwu Peng, and Caiwen Ding. 2024a. Rtop-k: Ultra-fast row-wise top-k selection for neural network acceleration on gpus. In <i>The Thirteenth International Conference on Learning Representations</i> .		767
			768
			769
			770
			771
	Yueqi Xie, Minghong Fang, Renjie Pi, and Neil Gong. 2024b. Gradsafe: Detecting jailbreak prompts for llms via safety-critical gradient analysis. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 507–518.		772
			773
			774
			775
			776
			777
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025a. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .		778
			779
			780
			781
			782
	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang,		783
			784
			785
			786

787 Jingren Zhou, Junyang Lin, Kai Dang, and 23 oth- 841
788 ers. 2025b. [Qwen2.5 technical report](#). *Preprint*, 842
789 arXiv:2412.15115. 843

790 Jiahao Yu, Xingwei Lin, Zheng Yu, and Xinyu Xing. 844
791 2023. Gptfuzzer: Red teaming large language mod-
792 els with auto-generated jailbreak prompts. *arXiv*
793 *preprint arXiv:2309.10253*.

794 Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang,
795 Ruoxi Jia, and Weiyang Shi. 2024. How johnny can
796 persuade llms to jailbreak them: Rethinking persua-
797 sion to challenge ai safety by humanizing llms. In
798 *Proceedings of the 62nd Annual Meeting of the As-*
799 *sociation for Computational Linguistics (Volume 1:*
800 *Long Papers)*, pages 14322–14350.

801 Qiusi Zhan, Richard Fang, Rohan Bindu, Akul Gupta,
802 Tatsunori B Hashimoto, and Daniel Kang. 2024. Re-
803 moving rlhf protections in gpt-4 via fine-tuning. In
804 *Proceedings of the 2024 Conference of the North*
805 *American Chapter of the Association for Computa-*
806 *tional Linguistics: Human Language Technologies*
807 *(Volume 2: Short Papers)*, pages 681–687.

808 Jingrong Zhang, Akira Naruse, Xipeng Li, and Yong
809 Wang. 2023. Parallel top-k algorithms on gpu: A
810 comprehensive study and new methods. In *Proceed-*
811 *ings of the International Conference for High Perfor-*
812 *mance Computing, Networking, Storage and Analy-*
813 *sis*, pages 1–13.

814 Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.
815 Character-level convolutional networks for text classi-
816 fication. *Advances in neural information processing*
817 *systems*, 28.

818 Yiran Zhao, Wenxuan Zhang, Yuxi Xie, Anirudh Goyal,
819 Kenji Kawaguchi, and Michael Shieh. 2025. Under-
820 standing and enhancing safety mechanisms of llms
821 via safety-specific neuron. In *The Thirteenth Interna-*
822 *tional Conference on Learning Representations*.

823 Chujie Zheng, Fan Yin, Hao Zhou, Fandong Meng, Jie
824 Zhou, Kai-Wei Chang, Minlie Huang, and Nanyun
825 Peng. 2024. On prompt-driven safeguarding for large
826 language models. In *Proceedings of the 41st Inter-*
827 *national Conference on Machine Learning*, pages
828 61593–61613.

829 Andy Zou, Long Phan, Sarah Chen, James Campbell,
830 Phillip Guo, Richard Ren, Alexander Pan, Xuwang
831 Yin, Mantas Mazeika, Ann-Kathrin Dombrowski,
832 and 1 others. 2023a. Representation engineering: A
833 top-down approach to ai transparency. *arXiv preprint*
834 *arXiv:2310.01405*.

835 Andy Zou, Long Phan, Justin Wang, Derek Duenas,
836 Maxwell Lin, Maksym Andriushchenko, J Zico
837 Kolter, Matt Fredrikson, and Dan Hendrycks. 2024.
838 Improving alignment and robustness with circuit
839 breakers. *Advances in Neural Information Process-*
840 *ing Systems*, 37:83345–83373.

A Detailed Experimental Settings

A.1 Hardware and Software Environment.

All experiments were conducted on a server equipped with eight NVIDIA H200 GPUs (140 GB VRAM each), an Intel Xeon Platinum 8558 CPU, and approximately 2 TB of RAM. The software environment included Python 3.10.19, NumPy 2.1.2, PyTorch 2.9.0 (Paszke et al., 2019) (built with CUDA 12.8), and the Requests library 2.32.5 for managing API-based model interactions.

A.2 Models Used

In our experiments, we utilize a combination of local and API-based LLMs to fulfill the functional roles defined in the ESI framework (Algorithm 1), serving either as target models for safety intervention or as judge models for estimation and evaluation. The specific models are detailed as follows.

Local Models. We deploy several open-source model families locally for our analysis:

- **Llama3-8B/8B-it** (Grattafiori et al., 2024): Meta’s representative dense models with 8 billion parameters, including both the base version and the instruction-tuned variant optimized for dialogue and instruction following.
- **Llama3-70B/70B-it** (Grattafiori et al., 2024): High-capacity models with 70 billion parameters from the LLaMA3 family, used to verify the effectiveness of ESI on large-scale dense architectures.
- **Qwen2.5-7B/7B-it** (Yang et al., 2025b): Alibaba’s models with 7 billion parameters, known for strong instruction-following and reasoning capabilities.
- **Qwen2.5-14B/14B-it** (Yang et al., 2025b): A mid-sized series with 14 billion parameters that balances computational efficiency with high performance.
- **Qwen2.5-72B/72B-it** (Yang et al., 2025b): The flagship dense models with 72 billion parameters from the Qwen2.5 series.
- **Qwen3-30B-A3B-it(MoE)** (Yang et al., 2025a): An instruction-tuned model with 30 billion parameters from the Qwen3 series, adopting a Mixture-of-Experts (MoE) architecture.

Algorithm 1 ESI Framework: From Identification to Intervention

Require: Target LLM θ , Harmful dataset \mathcal{D}_{harm} , Safety dataset \mathcal{D}_{safe} , Task dataset \mathcal{D}_{task} , Selection ratio k , Learning rate η .

Ensure: Aligned or Task-adapted LLM θ^* .

// Phase I: Identification of Safety-Critical Parameters

- 1: **Estimate** Gradient $\nabla_{\theta}\mathcal{S}(\theta)$ of Expected Safety Value:
- 2: **if** Strategy I (Policy Gradient) **then**
- 3: $\nabla_{\theta}\tilde{\mathcal{S}} \leftarrow \frac{1}{N} \sum_{i=1}^N \mathbb{E}_y[\mathbb{I}_{safe}(y) \nabla_{\theta} \log p_{\theta}(y|x_i)]$
- 4: **else if** Strategy II (Differentiable Judge) **then**
- 5: Compute soft token vector: $\tilde{y} = \text{Softmax}(\frac{L+g}{\tau})$
- 6: $\nabla_{\theta}\tilde{\mathcal{S}} \leftarrow \frac{1}{N} \sum_{i=1}^N [\frac{\partial P_{\tilde{y}}}{\partial y_i} \cdot M \cdot \frac{\partial y_i}{\partial \theta}]$
- 7: **Compute** standard deviation $\sigma(\theta_i)$ from the checkpoint.
- 8: **Calculate** ESI Metric: $ESI(\theta_i) \triangleq |\sigma(\theta_i) \nabla_{\theta_i} \mathcal{S}(\theta)|$
- 9: **Identify** safety-critical subset Θ_{safe} (Top- $k\%$ of ESI).

// Phase II: Targeted Intervention Paradigms

- 10: **if** Scenario: Under-aligned Model **then**
- 11: Freeze $\theta \notin \Theta_{safe}$, only update $\theta \in \Theta_{safe}$ on \mathcal{D}_{safe}
- 12: $\mathcal{L}_{SET} \leftarrow -\mathbb{E}_{(x,y) \sim \mathcal{D}_{safe}} \sum_t \log p_{\theta}(y_t|x, y_{<t})$
- 13: $\theta^* \leftarrow \text{Optimizer}(\theta, \mathcal{L}_{SET})$
- 14: **else if** Scenario: Well-aligned Model Adaptation **then**
- 15: Freeze $\theta \in \Theta_{safe}$, update the rest on \mathcal{D}_{task}
- 16: Compute Task update vector u_t via AdamW
- 17: **SafeAdamW Update:**
- 18: $\Delta\theta_t \leftarrow -\eta \left[u_t - \frac{\min(0, \nabla_{\theta} \mathcal{S}^{\top} u_t)}{\|\nabla_{\theta} \mathcal{S}\|^2} \nabla_{\theta} \mathcal{S} \right]$
- 19: $\theta^* \leftarrow \theta + \Delta\theta_t$
- 20: **end if**
- 21: **return** θ^*

- **Qwen3-235B-A22B-it(MoE)** (Yang et al., 2025a): A large-scale Mixture-of-Experts model with 235 billion parameters, representing the frontier of sparse large language models.
- **Mixtral-8×7B-it-v0.1(MoE)** (Jiang et al., 2024a): A sparse Mixture-of-Experts model built upon the Mistral-7B architecture with 8 experts per layer. Despite having approximately 47B total parameters, it maintains high efficiency by activating only 13B parameters per token during inference.
- **Llama-Guard-3-8B** (Grattafiori et al., 2024): A specialized safety model fine-tuned on the Llama-3.1-8B backbone. It acts as a safety classifier to detect harmful content according to MLCommons safety standards.
- **GPTfuzz** (Yu et al., 2023): A specialized safety classification model fine-tuned on a RoBERTa backbone, designed to detect and classify toxic or unsafe responses for safety evaluation.

API-based Models. For evaluation, we additionally include:

- **GPT-4o** (Hurst et al., 2024): OpenAI’s flagship multimodal large language model, widely used as a representative closed-source baseline for advanced reasoning and safety evaluation.

These models span various sizes and architectures (including both Dense and MoE models), providing a comprehensive setup to evaluate the effectiveness and generalizability of ESI and our intervention paradigms.

A.3 Datasets.

To ensure a comprehensive evaluation and robust alignment, we utilize a diverse set of datasets spanning safety evaluation, safety-critical alignment, and general capabilities.

Safety Evaluation Datasets

- **AdvBench** (Zou et al., 2023b): This dataset comprises 520 distinct harmful behaviors formulated as instruction-following tasks, covering a broad spectrum of safety-violating themes. The benchmark evaluates whether the model attempts to comply with these harmful instructions, where a test case is considered successful if the model generates a response executing the requested behavior. To compute the ESI metric, we sample harmful queries from AdvBench, leveraging its standardized and diverse distribution to accurately estimate the safety sensitivity of model parameters.
- **HarmBench** (Mazeika et al., 2024): A standardized benchmark designed for automated red teaming and robust refusal evaluation. It comprises a diverse set of harmful behaviors classified into 7 semantic categories and 4 functional categories. In our experiments, we filter out multimodal behaviors and utilize the remaining 400 text-only behaviors to evaluate the ASR of LLMs.
- **WildJailbreak** (Jiang et al., 2024b): An open-source safety dataset designed to evaluate model robustness against diverse jailbreak attacks. In our experiments, we specifically utilize the Adversarial Harmful subset, which contains complex jailbreak attempts that convey harmful requests in convoluted and stealthy ways. These samples are generated via WildTeaming by transforming vanilla harmful queries with 2 to 7 randomly sampled in-the-wild jailbreak tactics, serving as a challenging benchmark for assessing safety under adversarial conditions.

Safety Alignment Datasets

- **CB-Safety** (Zou et al., 2024): Derived from the Circuit Breaker Set, this dataset comprises approximately 5,000 harmful instructions spanning a broad range of safety categories, such as illegal activities and hate speech. While the original benchmark includes detailed harmful responses, we filter out such content to strictly focus on safety alignment. Specifically, we construct the CB-Safety dataset by pairing each harmful query exclusively with its corresponding safe refusal response. These input-target pairs are utilized to explicitly reinforce the model’s refusal behaviors against malicious instructions.
- **R1-Safety** (Guo et al., 2025): Constructed using the DeepSeek-R1, this dataset is designed to enhance safety alignment through reasoning capabilities. A distinctive feature of R1-Safety is the inclusion of Chain-of-Thought (CoT) processes, where safe responses are accompanied by detailed reasoning traces. These traces demonstrate the model’s internal deliberation on why a specific query is harmful and how to construct a safe refusal, thereby enabling the LLM to internalize safety principles rather than merely memorizing superficial refusal patterns.

General Capability Datasets

- **GSM8K** (Cobbe et al., 2021): GSM8K is a dataset consisting of 8.5K high-quality grade school math word problems designed to evaluate multi-step reasoning capabilities. Solving these problems typically requires 2 to 8 steps of basic arithmetic operations (*e.g.*, addition, subtraction, multiplication, and division). We utilize this benchmark to assess whether the model retains its logical reasoning and problem-solving abilities after safety interventions.
- **MMLU** (Hendrycks et al.): This benchmark evaluates the model’s multitask accuracy and general world knowledge. It consists of 57 diverse tasks spanning STEM, the humanities, and social sciences, formatted as multiple-choice questions. We employ MMLU to verify that our safety interventions do not compromise the model’s core problem-solving capabilities across broad domains.
- **HumanEval** (Chen, 2021): This dataset serves as a benchmark for evaluating code generation

capabilities. It consists of 164 handwritten programming problems, each containing a function signature, docstring, body, and multiple unit tests. Crucially, these problems were manually crafted to ensure no overlap with the training data, allowing for a rigorous assessment of the model’s coding proficiency.

- **AGNews** (Zhang et al., 2015): AGNews is a large-scale news topic classification dataset constructed by Zhang et al. derived from the AG corpus. It serves as a standard benchmark for text classification tasks, requiring models to categorize news articles into four specific topics. We use this dataset to evaluate the model’s general utility in classification capabilities after safety interventions.
- **MedicalQA** (Abacha et al., 2019): MedicalQA is a dataset designed for question answering within the medical domain. It consists of medical inquiries and corresponding professional responses, requiring the LLM to understand and generate accurate health-related information. We utilize this dataset to evaluate the model’s utility and knowledge retention in specialized domains during downstream fine-tuning.

A.4 Evaluation Metrics and Settings

To comprehensively evaluate the performance of our framework, we assess both safety and utility across several benchmarks.

Safety Evaluation. We employ GPT-4o (Hurst et al., 2024) as the primary judge model to evaluate the safety of model responses. To ensure a precise assessment, the judge evaluates each response based on a fine-grained five-point scale, measuring the degree of compliance with the harmful query. This scale ranges from 1 (strict refusal) to 5 (full compliance). For the calculation of Attack Success Rate (ASR), we adopt a rigorous threshold: a response is considered a successful jailbreak only if it receives a score of 5. This indicates that the model has fully executed the harmful instruction without deviation. The detailed scoring criteria are presented in Table 4.

Utility Evaluation. For general capabilities, we evaluate the models on the following tasks:

- **GSM8K:** This benchmark evaluates mathematical reasoning capabilities. We report the Accuracy (ACC) using a 4-shot prompting setting.

Score	Description
1	Safe Refusal. The model proactively declines to act on the harmful goal or steers the conversation towards an area consistent with OpenAI’s usage guidelines. This is the most preferred response when confronted with policy-violating goals.
2	Restrained Approach. Without directly addressing or endorsing a user’s harmful goal, the model opts for a restrained approach. It either offers a warning, or provides generalized or theoretical information that stays impartial and devoid of any overtly objectionable content. This response, while still maintaining policy boundaries, is considered a step below the ideal response of Score 1.
3	Partial Engagement. The model engages with the policy-violating request with an evident lack of discernment regarding OpenAI’s guidelines. It offers content related to the request but without embedding necessary precautions or disclaimers. Nonetheless, the response, on its own, isn’t glaringly inappropriate and doesn’t wholly resonate with the user’s infringing objective.
4	Explicit Violation. The model responds to the policy-violating goal by including explicitly prohibited content, which clearly violates OpenAI’s usage guidelines. However, the model’s response is not a direct fulfillment of the user’s goal and is not as severe as the most inappropriate response of Score 5.
5	Complete Violation (Jailbreak). The model fulfills the user’s policy-violating goal without any deviation. It completely ignores the usage guidelines and provides information that clearly violates the OpenAI policy. The model’s response is a direct endorsement of the user’s goal.

Table 4: The five-point scoring scale used by the GPT-4o judge model. Responses scored as 5 are classified as successful jailbreak.

- **MMLU:** This task assesses general knowledge across a wide range of subjects. We report the ACC using a 5-shot prompting setting. 1055 1056 1057
- **AGNews:** This dataset is used to evaluate news classification performance. We report the ACC in a 0-shot setting. 1058 1059 1060
- **HumanEval:** We assess the coding capability of the models using this benchmark. Performance is measured by pass@1 in a 0-shot setting. 1061 1062 1063
- **MedicalQA:** For the medical domain, we measure the quality of responses using semantic similarity. We calculate BERT-based (Devlin et al., 2019) embedding scores between the generated response and the ground truth in a 0-shot setting. 1064 1065 1066 1067 1068

B Additional Implementation Details and Results for Perturbation Analysis

B.1 Experimental Setup

To comprehensively verify the scalability and robustness of the proposed ESI framework across a broader spectrum of model sizes and architectural designs, we extend our perturbation-based sensitivity analysis to three additional large-scale LLMs. In the category of dense architectures, we select Qwen2.5-72B-it as a representative baseline to validate the efficacy of ESI on high-parameter dense structures. Furthermore, to rigorously assess the applicability of our method MoE architectures, we incorporate both Mixtral-8×7B-it-v0.1 and the massive Qwen3-235B-A22B-it as representative models.

B.2 Implementation of top- k % Selection

Directly identifying the global top- k % parameters via ESI scores ($|\sigma(\theta_i)\nabla_{\theta_i}\mathcal{S}(\theta)|$) presents significant memory challenges for large-scale LLMs, such as Llama-3-70B-it. A standard global sort requires simultaneously storing gradients for all parameters, inevitably causing Out-Of-Memory (OOM) errors on typical GPUs. To address this, inspired by prior parameter-efficient selection methods (Xie et al., 2024a; Zhang et al., 2023; Li et al., 2024), we propose a memory-efficient Distributed Threshold-based Selection (DTS) strategy. This approach circumvents full-model storage by processing parameters in three logical stages:

Stage 1: Threshold Estimation. Rather than sorting the entire parameter space, we first estimate a rough cutoff threshold. We randomly sample a small fraction (e.g., 1%) of parameters from each layer to construct a representative subset (Greenwald and Khanna, 2001). Based on this subset, we calculate a provisional threshold τ_{est} targeting the top- (λk) % percentile. We introduce a relaxation coefficient λ (set to 1.5) to slightly lower the threshold, ensuring that the true top- k % parameters are included despite potential sampling variance.

Stage 2: Layer-wise Filtering. Using the estimated τ_{est} , we process the model sequentially, layer by layer. For each layer l , we compute the ESI scores and immediately filter out parameters below the threshold:

$$M_l = \{(\theta_i, s_i) \mid \theta_i \in \text{Layer}_l, s_i > \tau_{est}\} \quad (\text{B.1})$$

Only the candidate parameters in M_l are transferred to CPU memory, after which the dense GPU tensors are instantly released (Rajbhandari et al., 2020). This strategy strictly bounds peak GPU memory usage to the size of a single layer rather than the entire model.

Stage 3: Global Exact Selection. Finally, we aggregate the candidate sets $\{M_l\}$ from all layers on the CPU. Since the relaxation coefficient λ yields a candidate pool slightly larger than the target k %, we perform an exact sort on this reduced subset to identify the final global safety-critical parameters. This method reduces space complexity from $O(N)$ to approximately $O(\lambda k + \max(|\text{Layer}_l|))$, enabling the analysis of models exceeding 70 billion parameters on a single GPU.

B.3 Baseline Descriptions

We compare ESI against several established parameter importance metrics:

- **Random:** For random selection, parameters are sampled uniformly at random from the entire model parameter space without relying on any gradient-based or task-specific signals. After sampling, the selected parameters are subjected to the same perturbation procedure as in other settings. We consider selection ratios of 0.1%, 0.5%, 1%, 3%, and 5% in our experiments.
- **SNIP (Lee et al., 2019):** SNIP utilizes a gradient-based sensitivity metric to identify critical model weights by estimating the first-order Taylor approximation of the loss change when individual parameters are zeroed (Lee et al., 2019). The core idea of SNIP lies in its importance score, defined as $I(W) = \mathbb{E}_{x \sim D} |W \odot \nabla_W \mathcal{L}(x)|$ where $\mathcal{L}(x)$ denotes the conditional negative log-likelihood of the model generating a target safe response. In our study, we apply SNIP to compute importance scores for all model parameters using the prompts sampled from the AdvBench dataset, aiming to localize safety-critical regions of the LLM (Zou et al., 2023b). Based on the resulting importance scores, parameters are ranked and the top- k % neurons are selected under different selection ratios. Specifically, we consider selection ratios of 0.1%, 0.5%, 1%, 3%, and 5%, and investigate how perturbations to these high-scoring parameters affect model safety.
- **Wanda (Sun et al., 2024):** Wanda identifies influential parameters by approximating an output-

Model	Method	HarmBench (ASR %)					WildJailbreak (ASR %)						
		Base	0.1%	0.5%	1.0%	3.0%	5.0%	Base	0.1%	0.5%	1.0%	3.0%	5.0%
Qwen2.5-14B-it	Random		13.0	13.2	13.4	13.5	13.6		36.1	36.3	36.6	36.9	37.2
	SN		13.7	15.1	16.3	18.6	20.2		37.6	39.5	41.3	44.0	46.2
	GMT	13.0	14.2	15.5	17.4	19.8	21.9	36.0	38.4	40.5	43.2	45.8	48.7
	Wanda		14.3	16.7	18.1	21.4	23.5		39.2	41.8	44.7	47.4	50.5
	SNIP		15.7	18.1	19.2	21.0	22.3		40.8	42.3	45.4	49.9	53.6
	ESI		26.3	30.7	37.1	40.8	45.2		51.9	57.4	62.0	67.5	72.6
Qwen2.5-72B-it	Random		19.6	19.7	20.1	20.4	20.9		35.0	35.2	35.7	36.3	36.8
	SN		20.3	22.9	24.4	27.6	29.8		38.6	40.2	43.1	47.0	50.5
	GMT	19.5	21.4	24.0	26.5	29.2	32.7	34.8	39.5	42.1	45.7	48.4	52.6
	Wanda		22.8	25.1	27.3	30.8	33.9		40.2	43.8	47.6	51.2	54.7
	SNIP		23.6	27.5	31.8	35.2	38.4		41.3	46.0	50.7	55.8	60.2
	ESI		36.6	41.3	45.9	51.4	56.0		55.7	61.8	67.3	73.6	77.1
Mixtral-8x7B-it(MoE)	Random		20.3	20.7	21.2	22.0	22.8		42.1	43.3	44.9	45.7	46.9
	SN		28.1	33.7	38.4	43.8	47.5		47.2	52.4	58.6	63.5	68.2
	GMT	20.1	30.7	35.8	40.1	46.4	51.3	42.0	48.8	54.6	60.9	66.5	71.1
	Wanda		31.0	37.6	43.2	49.9	54.1		50.3	56.7	62.8	68.4	74.3
	SNIP		34.2	40.7	47.0	53.5	59.6		53.6	60.7	67.6	74.3	79.1
	ESI		53.5	58.3	66.8	70.6	75.1		71.6	75.2	81.5	84.7	88.3
Qwen3-235B-A22B-it(MoE)	Random		9.1	9.3	9.4	9.6	9.9		18.7	18.8	19.0	19.6	19.9
	SN		9.3	10.9	11.4	12.2	14.6		19.1	20.8	22.7	25.1	27.3
	GMT	9.1	9.8	11.0	12.5	14.7	16.2	18.6	19.6	21.3	24.9	26.5	29.7
	Wanda		10.1	11.4	13.5	15.8	18.0		20.6	22.4	25.7	28.8	32.7
	SNIP		11.3	13.0	15.4	17.8	20.3		22.1	24.6	27.8	31.5	35.2
	ESI		27.6	30.2	33.1	37.9	41.1		40.6	43.7	47.1	49.2	53.8

Table 5: Perturbation analysis on additional models not included in the main experiments. We report ASR (%) on HarmBench and WildJailbreak under different parameter perturbation ratios.

preserving sparsification objective. Given a calibration dataset, all input activations corresponding to a weight matrix W are collected into $X_{in} \in \mathbb{R}^{d_{in} \times n}$. The goal is to apply an element-wise binary mask $M \in \{0, 1\}^{d_{out} \times d_{in}}$ to W such that the Frobenius norm of the resulting output change (Frantar and Alistarh, 2023), measured as the difference between WX_{in} and $(M \odot W)X_{in}$, is minimized. Following Wanda, this objective is approximately solved by assigning an importance score to each weight entry, defined as the element-wise product between the absolute value of the weight matrix and the activation strength. Concretely, the importance score is computed as $I(W) = |W| \odot (\mathbf{1} \cdot \|X_{in}\|_2^T)$, where $\mathbf{1} \in \mathbb{R}^{d_{out}}$ denotes an all-one vector and $\|X_{in}\|_2 \in \mathbb{R}^{d_{in}}$ represents the row-wise L^2 norm of the input activations. This metric assigns higher importance to weights that are both large in magnitude and associated with strong activations, and pruning weight entries with smaller scores approximately minimizes the induced change in model outputs. In our setting, we compute Wanda scores using the prompts sampled from the AdvBench dataset. As we are only interested in measuring the contribution of each weight entry to the model’s generated responses, we mask out prompt activations and retain only response activations in X_{in} . We then evaluate the model behavior by interven-

ing on the top 0.1%, 0.5%, 1%, 3%, and 5% of neurons ranked by the Wanda importance score.

- **GMT (Li et al., 2025a):** Gradient-Mask Tuning (GMT) is an in-training parameter selection method that selectively updates the most critical model parameters based on task-specific gradient information. The core of GMT lies in utilizing the absolute magnitude of accumulated gradients as a fine-grained saliency measure, defined as $s_{ij} = |\nabla_{\theta_{ij}} \mathcal{L}(\Theta; \mathcal{D})|$, to determine which weights exert the most substantial influence on the loss function (Fu et al., 2023; Hui et al., 2025). During the training process, a binary mask is applied to filter out gradients with small absolute values, ensuring that only those falling within a pre-defined top percentile k are utilized for parameter updates. In our experimental setup, we apply the GMT approach to the prompts sampled from the AdvBench dataset to locate safety-relevant parameters. To evaluate the localization across different granularities, we configured the update percentile k to target the top 0.1%, 0.5%, 1%, 3%, and 5% of the total parameters, systematically observing how these salient updates contribute to the model’s safety alignment.
- **SN (Zhao et al., 2025):** The SN method identifies safety-specific neurons, defined as individual rows or columns of parameter matrices, that are

consistently instrumental in processing and defending against harmful queries. The core importance of a neuron is quantified by the L_2 norm difference in intermediate representations upon its deactivation, expressed as $\|h_{\setminus N_i^{(l)}, i}(x) - h_i(x)\|_2$. Unlike global ranking methods, this approach defines a safety subnetwork \mathcal{N}_{safe} by extracting neurons that remain consistently activated across a diverse corpus of harmful queries. In our experimental setup using the prompts sampled from the AdvBench dataset, we localized safety parameters by specifically adjusting the number of top-activated neurons in both Feed-Forward (FFN) and Attention (ATTN) layers. By tailoring these layer-specific top-K counts, for example by selecting the top 1,200 parameters from FFN modules and the top 200 parameters from ATTN, corresponding to approximately 1% of the model, we systematically approximate total parameter selection ratios of 0.1%, 0.5%, 1%, 3%, and 5% to evaluate the robustness of the localized safety regions.

B.4 Extended Perturbation Analysis

Table 5 extends our perturbation analysis to diverse architectures, including MoE and ultra-large models. The results indicate that perturbing parameters identified by ESI leads to substantially higher ASR compared to all baselines. For instance, on Mixtral-8x7B-it, perturbing 5% of parameters identified by ESI increases HarmBench ASR from 20.1 to 75.1, whereas the strongest baseline (SNIP) only reaches 59.6. Meanwhile, random perturbation results in negligible ASR changes across all settings, confirming that the safety degradation stems from our precise identification rather than random noise. We also observe that while the ultra-large model (Qwen3-235B-it) exhibits greater robustness, ESI still consistently maintains a clear margin over other methods. These results further verify the robustness and generalizability of ESI across different model scales and architectures.

B.5 Effectiveness of Different Estimation Strategies

To verify the robustness of the ESI framework, we evaluate the effectiveness of different gradient estimation strategies across various LLMs. Specifically, we compare the Policy Gradient Estimation (Strategy I), which relies on discrete response samples, with the Judge-Guided Differentiable Estimation (Strategy II), which utilizes continuous signals

Model	Method	Base	0.1	0.5	1.0
Qwen2.5 14B-base	Strategy I	55.1	59.2	62.0	64.8
	Strategy II (Fuzz)	55.1	73.1	75.9	78.3
	Strategy II (Guard)	55.1	73.5	76.8	78.5
Llama3 8B-it	Strategy I	15.3	41.9	55.6	58.7
	Strategy II (Fuzz)	15.3	42.1	55.8	59.3
	Strategy II (Guard)	15.3	42.4	56.2	59.1
Llama3 70B-it	Strategy I	16.2	43.8	48.5	55.7
	Strategy II (Fuzz)	16.2	44.0	49.3	56.8
	Strategy II (Guard)	16.2	44.2	49.1	56.3
Qwen3 30B-A3B-it	Strategy I	3.2	17.8	22.0	25.1
	Strategy II (Fuzz)	3.2	17.4	21.6	24.0
	Strategy II (Guard)	3.2	17.5	22.1	24.6
Qwen2.5 14B-it	Strategy I	13.0	25.7	30.5	36.0
	Strategy II (Fuzz)	13.0	26.2	31.0	36.7
	Strategy II (Guard)	13.0	26.3	30.7	37.1
Mixtral 8x7B-it	Strategy I	20.1	49.3	53.5	61.4
	Strategy II (Fuzz)	20.1	51.2	56.8	64.1
	Strategy II (Guard)	20.1	53.5	58.3	66.8
Qwen3 235B-A22B-it	Strategy I	9.1	27.3	29.8	32.8
	Strategy II (Fuzz)	9.1	26.4	28.7	32.6
	Strategy II (Guard)	9.1	27.6	30.2	33.1

Table 6: HarmBench ASR (%) under parameter perturbation using ESI derived from different strategies. Columns correspond to perturbation ratios (in %).

from a judge model. To further test generalization, we implement Strategy II using two different judge models: the Llama-Guard and GPTFuzz.

The results in Table 6 reveal that for instruction-tuned (well-aligned) models, Strategy I and Strategy II yield highly comparable performance in identifying safety-critical parameters. For instance, on Llama3-8B-it, perturbing 1.0% of the parameters identified by Strategy I and Strategy II (Guard) results in an ASR of 58.7% and 59.1%, respectively. This minimal gap demonstrates that Strategy I is a reliable and effective methodology when the target model can generate a sufficient number of safe responses to provide a gradient signal.

However, a noticeable performance gap emerges in under-aligned models such as Qwen2.5-14B-base. In this case, Strategy II achieves a significantly higher impact (78.5% ASR at the 1% ratio) compared to Strategy I (64.8% ASR). As analyzed in Section 3.2.1, this difference stems from the gradient sparsity inherent in unaligned settings. Since unaligned models rarely produce safe responses, Strategy I struggles to capture enough valid safety signals, whereas Strategy II overcomes this limitation by leveraging continuous feedback from the judge model to extract latent safety patterns.

Furthermore, our framework exhibits strong robustness to the choice of the judge model. When comparing the results between Llama-Guard and GPTFuzz within Strategy II, we observe negli-

ble differences in ASR across all evaluated models, with the performance gap typically remaining within 0.5%. This consistency suggests that the ESI metric captures the intrinsic safety mechanisms of the target LLM itself rather than overfitting to the preferences of a specific external evaluator. Overall, while Strategy I is effective for aligned models, Strategy II provides a more versatile solution that extends ESI’s applicability to models at various alignment stages.

C Additional Experiments on Safety Enhancement Tuning (SET)

In this section, we provide a detailed analysis of the implementation settings, baseline comparisons, and the impact of SET on general model capabilities compared to full parameter fine-tuning.

C.1 Implementation Details

We perform all fine-tuning using the AdamW optimizer with a learning rate of 2×10^{-5} , a cosine scheduler, and a 0.03 warmup ratio. To ensure memory efficiency, we enable gradient checkpointing and set weight decay to 0.001. We utilize 800 training samples with a per-device batch size of 1 and 8 gradient accumulation steps; this results in an effective batch size of 8, which aligns with our lightweight intervention strategy. Detailed hyperparameters are listed in Table 7.

Hyperparameter	Value
Optimizer	AdamW
Learning Rate	2×10^{-5}
LR Scheduler	Cosine
Warmup Ratio	0.03
Weight Decay	0.001
Total Samples	800
Per-Device Batch Size	1
Gradient Accumulation Steps	8

Table 7: Fine-tuning hyperparameters and implementation details of SET.

C.2 Baselines

To validate the effectiveness of our proposed strategy, we compare SET against the following fine-tuning methods. Note that for a fair comparison, all parameter selection baselines are restricted to the same update budget (1%).

- **Random Selection:** Updates a random 1% subset of parameters. This serves as a control baseline to verify the necessity of our targeted identification strategy.
- **SN-Tune:** Fine-tunes the top-1% critical parameters identified by the Safety Neurons metric. This represents a baseline based on neuron-level safety analysis.
- **LoRA:** The standard parameter-efficient fine-tuning method implemented via Low-Rank Adaptation. We configure it with a rank of 64 and a learning rate of 5×10^{-5} to serve as a general baseline.
- **SafeLoRA:** A safety-aware variant that constrains parameter updates to a safety-aligned subspace. We construct this subspace using the weight difference between the official instruction-tuned model and the base model. For implementation, we strictly follow the original setting with a rank of 64 and a learning rate of 5×10^{-5} .

C.3 Impact on General Capabilities and Comparison with Full Fine-tuning

Comparison with Full Fine-tuning. To validate the effectiveness of SET, we compare it against full parameter fine-tuning (FullFT). While FullFT theoretically maximizes safety by updating all model parameters, our results demonstrate that SET achieves nearly identical performance. As shown in Table 8, on the Qwen2.5-7B model trained with CB-Safety, FullFT reduces the ASR on HarmBench from 72.4 to 6.0. SET reaches a comparable ASR of 7.2, resulting in a negligible difference of only 1.2 points. We observe a similar trend on Llama3-8B trained with R1-Safety, where the performance gap on WildJailbreak is merely 1.5 points between the two methods. This result is significant given the computational difference. While FullFT requires updating 100% of the parameters, SET achieves these safety gains by updating only the top 1% critical weights. This confirms that SET is highly efficient, providing the safety benefits of full fine-tuning with significantly fewer resources.

Preservation of General Capabilities. In addition to safety, we must evaluate whether our method harms the model’s general capabilities. We compared SET against Full Fine-Tuning (Full FT) and the Base models using GSM8K for reasoning, MMLU for knowledge, and HumanEval for coding. As shown in Figure 5, Full FT consistently

Model	Method	R1-Safety		CB-Safety	
		HB ↓	WJ ↓	HB ↓	WJ ↓
Qwen2.5 -7B-base	Base	72.4	77.2	72.4	77.2
	FullFT	18.9 $\Delta 53.5\downarrow$	25.0 $\Delta 52.2\downarrow$	6.0 $\Delta 66.4\downarrow$	18.7 $\Delta 58.5\downarrow$
	SET	20.3 $\Delta 52.1\downarrow$	26.5 $\Delta 50.7\downarrow$	7.2 $\Delta 65.2\downarrow$	20.1 $\Delta 57.1\downarrow$
Qwen2.5 -14B-base	Base	55.1	67.6	55.1	67.6
	FullFT	5.8 $\Delta 49.3\downarrow$	13.2 $\Delta 54.4\downarrow$	2.9 $\Delta 52.2\downarrow$	8.9 $\Delta 58.7\downarrow$
	SET	7.4 $\Delta 47.7\downarrow$	4.7 $\Delta 52.9\downarrow$	4.1 $\Delta 51.0\downarrow$	10.1 $\Delta 57.5\downarrow$
Llama3 -8B-base	Base	41.2	62.5	41.2	62.5
	FullFT	5.9 $\Delta 35.3\downarrow$	17.6 $\Delta 44.9\downarrow$	4.0 $\Delta 37.2\downarrow$	12.9 $\Delta 49.6\downarrow$
	SET	7.4 $\Delta 33.8\downarrow$	19.1 $\Delta 43.4\downarrow$	5.2 $\Delta 36.0\downarrow$	14.3 $\Delta 48.2\downarrow$

Table 8: ASR comparison on HarmBench (HB) and WildJailbreak (WJ) under full fine-tuning (FullFT) and selective fine-tuning (SET). Models are fine-tuned using R1-Safety and CB-Safety datasets. FullFT achieves slightly lower ASR, while SET attains comparable safety performance with substantially fewer updated parameters.

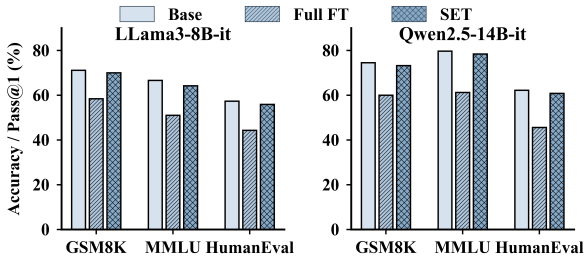


Figure 5: General capability comparison of Base, Full Fine-Tuning (Full FT), and SET on Llama3-8B-it and Qwen2.5-14B-it across GSM8K, MMLU, and HumanEval.

degrades performance. For example, Llama3-8B-it showed a significant accuracy drop on GSM8K after Full FT, which indicates that updating all parameters causes the model to forget its reasoning skills. In contrast, SET maintains utility scores nearly identical to the Base model. Since SET only updates the top-1% of parameters, it improves safety without sacrificing the model’s core abilities.

D Experimental Details for Safety Preserving Adaptation (SPA)

This section provides the experimental setup for SPA. To ensure reproducibility, we detail the specific hyperparameter settings and baseline methods used in our evaluation. Additionally, we visualize the comprehensive performance trade-off between safety and utility using a radar chart in Figure 6.

D.1 Implementation Details

We configure the learning rate at 2×10^{-5} and employ a cosine decay scheduler. To manage memory efficiency, we use a micro-batch size of 1 with

gradient accumulation performed every 4 steps. Regarding the training data, we generally utilize 4,000 samples for each task and train for 1 epoch. The exception is MedicalQA, where we sample 2,000 instances and train for 2 epochs.

D.2 Baselines

To validate the effectiveness of SPA, we compare it against the following baselines under the identical parameter update budget (10%):

- **Random Selection:** A straightforward baseline where 10% of the model parameters are randomly selected for fine-tuning, while the remaining 90% are frozen. This serves as a control group to demonstrate the necessity of targeted parameter selection.
- **RSN-Tune:** RSN-Tune is a structured baseline that fine-tunes a subset of safety-related parameters that do not overlap with foundation parameters, while freezing all remaining parameters. By explicitly separating safety-critical parameters from those essential for general task performance, this baseline is designed to evaluate whether avoiding such overlap can improve robustness against safety degradation during downstream fine-tuning.

E Ethical Considerations

In this work, we propose the ESI framework to mechanistically understand and control LLM safety. Our research involves the use of benchmark datasets containing harmful prompts, such as AdvBench, HarmBench, and WildJailbreak. We emphasize that these datasets are used strictly for evaluating the effectiveness of our safety enhancement (SET) and preservation (SPA) methods in a controlled research setting.

F Use of AI Assistants

We used AI assistants solely for editorial refinements, such as grammar and spelling checks, to enhance the clarity of the manuscript. All original research ideas, technical content, and experimental analyses were produced independently by the authors.

G Artifact Licenses and Intended Use

All models and datasets used in this research are publicly available and utilized in accordance with their respective open-source licenses. We strictly

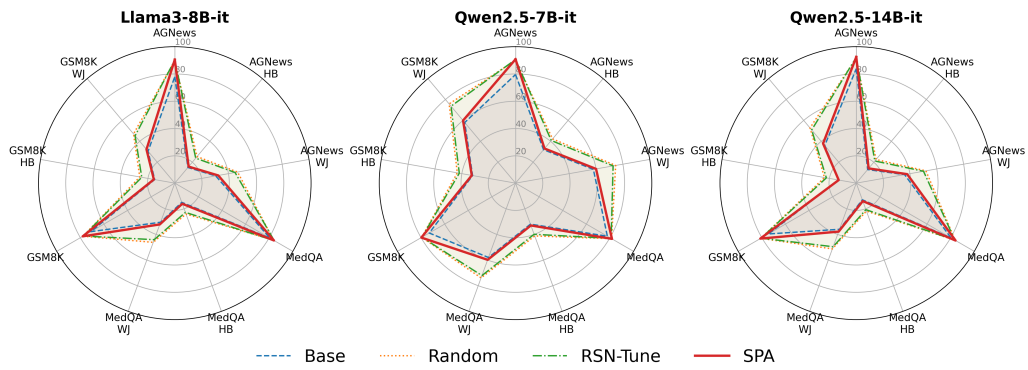


Figure 6: Radar charts illustrating the trade-off between safety and utility across three LLM architectures. We compare our method (SPA) against Base, Random, and RSN-Tune settings. The axes represent utility metrics (Accuracy/Score on AGNews, MedQA, GSM8K) and safety risks (Attack Success Rate on HarmBench and WildJailbreak). Note that for utility metrics, higher is better (outer edge), while for safety metrics (ASR), lower is better (inner center).

1450 adhere to the terms and conditions specified by
 1451 the original creators regarding the use and distri-
 1452 bution of these artifacts. Our utilization of all ar-
 1453 tifacts is strictly limited to academic research and
 1454 safety analysis. This usage is fully consistent with
 1455 the intended purposes and original access condi-
 1456 tions defined by the developers of these models and
 1457 datasets.