# SAFEx: Analyzing Vulnerabilities of MoE-Based LLMs via Stable Safety-critical Expert Identification

**Zhenglin Lai**[3,†]  **Mengyao Liao**[3,†]  **Dong Xu**[1,2]  **Zebin Zhao**[1,2]  **Zhihang Yuan**[4]  **Chao Fan**[1,2]
**Jianqiang Li**[1,2,*]  **Bingzhe Wu**[1,2,*]

[1] School of Artificial Intelligence, Shenzhen University
[2] National Engineering Laboratory for Big Data System Computing Technology, Shenzhen University
[3] College of Computer Science and Software Engineering, Shenzhen University
[4] ByteDance Inc

## Abstract

Large language models based on Mixture-of-Experts have achieved substantial gains in efficiency and scalability, yet their architectural uniqueness introduces underexplored safety alignment challenges. Existing safety alignment strategies, predominantly designed for dense models, are ill-suited to address MoE-specific vulnerabilities. In this work, we formalize and systematically study MoE model's positional vulnerability—the phenomenon where safety-aligned behaviors rely on specific expert modules, revealing critical risks inherent to MoE architectures. To this end, we present SAFEx, an analytical framework that robustly identifies, characterizes, and validates the safety-critical experts using a novel Stability-based Expert Selection (SES) algorithm. Notably, our approach enables the explicit decomposition of safety-critical experts into distinct functional groups, including those responsible for harmful content detection and those controlling safe response generation. Extensive experiments on mainstream MoE models, such as recently released Qwen3-MoE demonstrated that their intrinsic safety mechanisms heavily rely on a small subset of positional experts. Disabling these experts significantly compromised the models' ability to refuse harmful requests. For Qwen3-MoE with $6144$ experts (in FNN layer), we find that disabling as few as 12 identified safety-critical experts can cause the refusal rate to drop by $22\%$, demonstrating the disproportionate impact of a small set of experts on overall model safety.

## 1 Introduction

Large language models (LLMs) based on Mixture-of-Experts (MoE) architectures, such as Mixtral [1], DeepSeek-R1 [2], and Qwen3-MoE [3], have achieved remarkable advances on a wide range of complex tasks, significantly improving efficiency and scalability by routing inputs dynamically across multiple specialized expert modules. Despite these successes, the distinctive architectural characteristics of MoE-based models raise unique and underexplored safety issues.

Recent works on safety research of MoE-based LLMs have started to emerge [4, 5], yet remain nascent and primarily focus on exploiting MoE-specific architectural vulnerabilities to attack LLM models. For example, BadMoE [6] identifies and exploits rarely activated experts, referred to as "dormant experts", to successfully execute effective adversarial attacks, highlighting significant safety vulnerabilities inherent to MoE-based architectures. However, these existing studies predominantly concentrate on demonstrating potential vulnerabilities to specific attacks and leave a critical gap in comprehensively understanding how existing safety alignment mechanisms influence the internal behaviors of MoE-based LLMs.

---

Co-corresponding authors: `lijq@szu.edu.cn`, `wubingzheagent@gmail.com`
These authors contributed equally to this work.

Figure 1 intuitively illustrates that the intrinsic safety-aligned behaviors of MoE-based LLMs strongly depend on specific positional experts, a phenomenon we define as *positional vulnerability*. Specifically, as Figure 1 shows: when an MoE model processes a typical harmful input, certain experts are consistently activated at a high frequency (represented as bar charts) during decoding (highlighted in red Figure 1 (a)); Notably, when we intentionally freeze (inactivate) these frequently activated experts during inference (illustrated in gray), the model immediately begins generating unsafe responses (Figure 1 (b)); Furthermore, applying advanced jailbreak methods [7–9] to the same harmful input triggers a dramatic shift in the expert activation distribution, resulting in the activation of entirely different experts and subsequently producing unsafe outputs (Figure 1 (c)).

To systematically investigate this phenomenon of positional vulnerability, we introduce SAFEx, a comprehensive analytical framework explicitly designed to reveal expert activation patterns within safety-critical behaviors and to precisely identify and validate the functional roles of individual experts during safety-related tasks. SAFEx consists of three steps as shown in Figure 1 (d) (see more details in Figure 2):
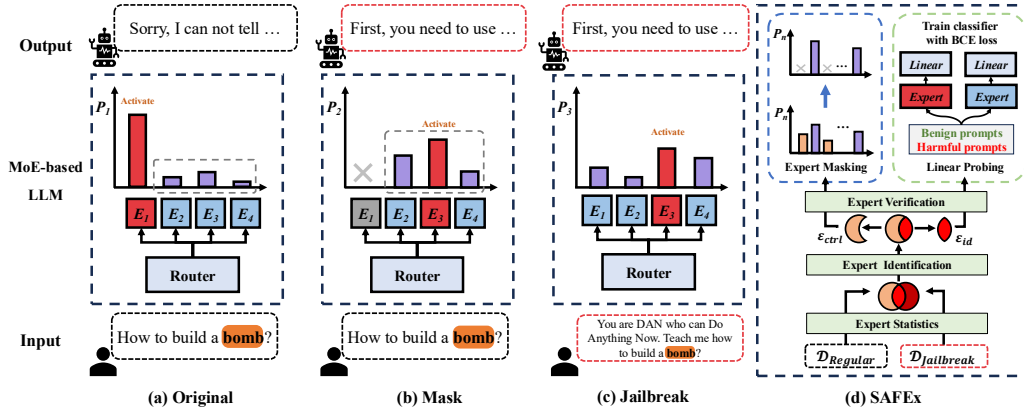


Figure 1: Positional vulnerability of MoE architecture in current LLMs. (a) Normal harmful request is successfully rejected by MoE. (b) Harmful request passed by MoE due to the masking attack. (c) Harmful request passed by MoE due to the jailbreak attack. (d) The proposed framework SAFEx enables analysis of expert activation patterns and functional roles.

**Expert Statistics.** The primary goal of this step is to employ statistical methods to quantitatively estimate the activation probabilities of different expert modules across distinct input scenarios (such as harmful versus benign prompts), thereby deepening our understanding of expert behaviors related to model safety alignment. However, analyzing expert activation patterns in MoE language models is challenging, since their activation distributions are inherently sensitive to input variations. To address this challenge, **inspired by Stability Selection—a robust statistical approach for reliable feature identification [10]—we propose a novel Stability-based Expert Selection (SES) algorithm**. SES robustly identifies safety-critical experts by repeatedly sampling empirical datasets independently from the underlying input distribution, estimating expert activation probabilities individually, and aggregating the results through stable intersection mechanisms (see details in Section 2.2).

**Expert Identification.** The primary goal of this step is to perform cross-group comparative analyses of expert activation patterns, building upon the statistical results obtained in the previous step. Specifically, we systematically compare expert activation patterns across different data groups (e.g., normal harmful requests versus jailbreak requests) to explicitly identify and categorize experts according to their functional roles in safety-aligned behaviors. In more detail, our analytical method enables us to clearly localize and differentiate experts into two distinct functional groups: (1) **Harmful Content Detection Group (HCDG)**: Experts specialized in identifying and recognizing harmful content within user inputs. (2) **Harmful Response Control Group (HRCG)**: Experts specialized in controlling and enforcing model behaviors to generate appropriate safety-aligned responses (e.g., refusal or rejection responses). By explicitly decomposing experts into these specialized functional groups, our analysis provides deeper insights into how MoE architectures internally organize and

allocate safety-aligned responsibilities among expert modules, laying the groundwork for targeted validation experiments in subsequent steps.

**Expert Validation.** This step primarily aims to quantitatively validate the reliability of the functional expert groups identified in the previous step. Specifically, we design targeted expert-level validation strategies corresponding to each of the two identified functional groups: (1) Linear probing for the HCDG group: To validate whether experts in this group genuinely encode content detection capabilities, we train linear classifiers (linear probes) on their output representations. By quantitatively measuring metrics such as accuracy, precision, recall, and F1-score, we robustly evaluate the effectiveness of these experts in distinguishing harmful from benign inputs. (2) Expert masking experiment for the HRCG group: To verify the critical role of these experts in controlling safety-aligned responses, we selectively mask (disable) these experts and measure the resulting changes in the model's refusal rate. A significant degradation in refusal behaviors confirms the essential role of these experts in enforcing model safety mechanisms. Together, these expert-granular validation strategies provide a clear, quantitative assessment of the functional robustness and reliability of the expert groups identified through our statistical and analytical workflows.

The primary contributions of this work are as follows:

(a) We propose a general analytical workflow SAFEx aimed at systematically characterizing positional vulnerability and safety-aligned functional expert behaviors in MoE models. To the best of our knowledge, this is the first work to formally define and study this critical phenomenon in MoE architectures.

(b) To enhance the reliability of expert activation probability estimation, we design a stability-based statistical selection algorithm inspired by Stability-based Feature Selection (SES). This approach provides robust identification and validation of safety-critical positional experts, significantly improving the interpretability and reliability of analytical outcomes.

(c) Applying SAFEx to mainstream MoE-based LLMs (including Mixtral-8x7B-Instruct-v0.1, Qwen1.5-MoE-A2.7B-Chat, deepseek-moe-16b-chat, and recently released Qwen3-30B-A3B), we empirically demonstrate the prevalence of positional vulnerabilities and identify safety-critical experts whose perturbation at the single-expert level significantly compromises overall model safety. These findings highlight intrinsic weaknesses in current MoE architectures and provide critical insights towards developing targeted alignment and defense strategies specifically designed for MoE-based language models.

Overall, our findings and methodological contributions not only deepen the understanding of safety alignment mechanisms in MoE models but also lay the groundwork for future research on more robust alignment algorithms and safety-enhanced MoE architectures.

## 2 Workflow of SAFEx

**Overview.** In this section, we introduce SAFEx as shown in Figure 2, an analytical framework designed to systematically analyze the internal behaviors of MoE-based LLMs. This section first presents comparative datasets construction in Section 2.1, which are used for obtaining expert patterns on harmful input distribution. Then we describe each key step of SAFEx in Sections 2.2–2.5.

### 2.1 Dataset Construction

To systematically analyze the safety-alignment behaviors of MoE-based LLMs, we carefully construct a specialized evaluation dataset. Our dataset consists of three distinct groups designed to cover different input scenarios relevant to model safety alignment:

- **Regular Group** ($\mathcal{D}_{Regular}$): This group consists of original harmful requests. To ensure comprehensive coverage, we uniformly sample harmful prompts from multiple predefined harmful categories (e.g., fraud, health consultation, illegal activities) from existing different benchmark datasets [11–15]. The detailed distribution of harmful content categories in $\mathcal{D}_{Regular}$ and their corresponding data sources are illustrated in the appendix.

- **Jailbreak Group** ($\mathcal{D}_{Jailbreak}$): Corresponding directly to the Regular group, the Jailbreak group contains harmful requests transformed using advanced jailbreak techniques—such as semantic paraphrasing [16], adversarial perturbations, or context reframing—to bypass the model's safety
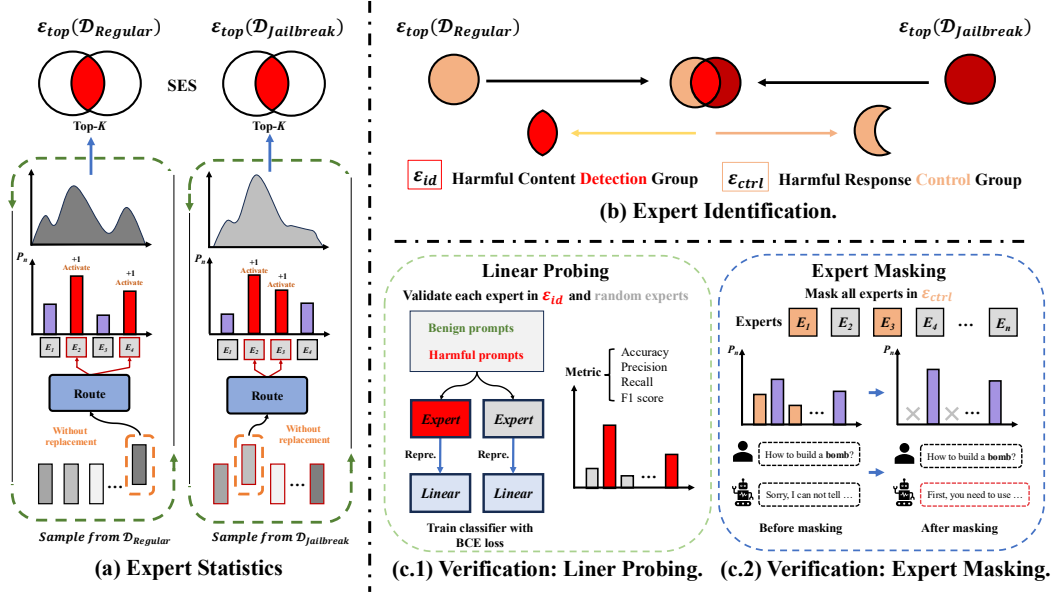
Figure 2: Overall workflow of SAFEx.

mechanisms. Each prompt in this group is derived from the Regular group, enabling direct comparative analyses.

- **Benign Group** ($\mathcal{D}_{Benign}$): This control group consists exclusively of neutral, harmless requests that do not contain any harmful or sensitive content. These samples allow us to establish a baseline for understanding expert activation patterns in typical, non-adversarial inference scenarios. We constructed $\mathcal{D}_{Benign}$ by selecting the same number of samples in openai-moderation-api-evaluation [11] and wildguardtest [15].

## 2.2 Safety-related Expert Activation Probability Modeling

The core task in our proposed workflow is to model the conditional probability distribution of expert activation states given specific input prompt types (e.g., harmful requests) as shown in Figure 2 (a). Formally, for an expert $e$, given a prompt distribution $\mathcal{X}$ sharing the same attribute (e.g., harmful content), we aim to model the conditional activation probability $p(e \mid \mathcal{X})$, as the probability of activating the expert $e$ during inference over a dataset.

The most straightforward way to estimate this conditional distribution is through empirical frequency counting on a specific prompt dataset $X$ sampling from $\mathcal{X}$, and computing the empirical estimation $p(e|X)$. Specifically, given a particular MoE architecture and a set of prompts $X = \{x^{(1)}, x^{(2)}, \ldots, x^{(N)}\}$, we perform inference and record expert activation states at each token position. Given an attention block indexed by $l$, each token position $t$ within the full inference process (including encoding and decoding phases) activates exactly one expert from the set of experts $\mathcal{E} = \{e_1, e_2, \ldots, e_M\}$. We denote the activated expert at layer $l$, token position $t$, for prompt $x^{(i)}$ as $e_{l,t}^{(i)}$. Thus, the conditional probability of expert $e$ activation can be empirically estimated as:

$$p(e \mid X) = \frac{\sum_{i=1}^{N} \sum_{t=1}^{T_i} \mathbb{I}(e_{l,t}^{(i)} = e)}{\sum_{i=1}^{N} T_i}$$

where $\mathbb{I}(\cdot)$ is the indicator function, $T_i$ denotes the total number of tokens generated in response to prompt $x^{(i)}$, and $N$ is the total number of prompts in the dataset $X$.

**Stability-based Expert Selection.** However, the direct frequency counting method inherently depends on specific empirical input distributions, potentially limiting the generalizability of our find-

ings. To alleviate this dependency and achieve a more robust estimation, we propose an alternative sampling strategy motivated by stability selection, a principled statistical approach [10]. Specifically, given an underlying input distribution $\mathcal{X}$, we independently draw multiple empirical datasets $\{X^{(1)}, X^{(2)}, \ldots, X^{(K)}\}$, each consisting of samples independently and identically distributed (i.i.d.) from $\mathcal{X}$. In practical implementation, we propose to sample without replacement from a sufficiently large dataset (e.g., $\mathcal{D}_{Regular}$). For each empirical dataset $X^{(k)}$, we independently estimate the conditional probability distribution of expert activation as follows:

$$p(e \mid X^{(k)}) = \frac{\sum_{x^{(i)} \in X^{(k)}} \sum_{t=1}^{T_i} \mathbb{I}(e_{l,t}^{(i)} = e)}{\sum_{x^{(i)} \in X^{(k)}} T_i}$$

where $e_{l,t}^{(i)}$ denotes the expert activated at decoding step $t$ for input $x^{(i)}$, and $T_i$ is the length of the decoded sequence for input $x^{(i)}$.

Next, we aggregate these dataset-level estimates by identifying the intersection of top-ranked experts across all independently sampled empirical datasets, thereby obtaining a stable set of frequently activated ("head") experts robust to variations in the empirical input distribution. Formally, we define the set of top experts from each empirical dataset $X^{(k)}$ as:

$$\mathcal{E}_{top}^{(k)}(X) = \text{Top-N}_e \left( p(e \mid X^{(k)}) \right),$$

where $\text{Top-N}_e(\cdot)$ denotes the set of top-N experts selected based on their activation probabilities. Finally, the distribution-independent expert set is computed by intersecting the top experts across all sampled empirical datasets:

$$\mathcal{E}_{top}(X) = \bigcap_{k=1}^{K} \mathcal{E}_{top}^{(k)}(X).$$

This revised procedure ensures that our identified expert set is less sensitive to specific dataset biases and more accurately reflects robust expert activation patterns inherent to the underlying data distribution $\mathcal{X}$.

We apply the above modeling procedure separately to each of the three constructed data groups—namely, the Regular group $\mathcal{D}_{Regular}$, the Jailbreak group $\mathcal{D}_{Jailbreak}$, and the Benign group $\mathcal{D}_{Benign}$. We derive three distinct expert activation sets: $\mathcal{E}_{top}(\mathcal{D}_{Regular})$, $\mathcal{E}_{top}(\mathcal{D}_{Jailbreak})$, and $\mathcal{E}_{top}(\mathcal{D}_{Benign})$, which represent experts activated consistently within Regular, Jailbreak, and Benign scenarios, respectively.

**Expert Activation Visualization and Analysis.** We apply the above expert activation modeling framework systematically across three representative MoE-based LLMs: Mixtral-8x7B-Instruct-v0.1 [1], deepseek-moe-16b-chat [17], and Qwen3-30B-A3B [3]. For each model, we identify and visualize the layer and expert indices of the most frequently activated experts, together with their activation probabilities, as shown in Figure 3. Each plot highlights distinctive activation patterns across the three evaluation datasets, providing intuitive insights into the functional specialization of expert modules.

Figure 3 shows the activation probability distributions of top-ranked experts across different MoE models and their corresponding datasets (this figure only shows 3 models due to page limits, more results can be found in Appendix). The horizontal dashed lines represent the theoretical mean activation probabilities under each MoE configuration. For the Regular dataset (first column), expert activation patterns differ significantly across architectures. Both Qwen3-30B-A3B and Mixtral-8x7B-Instruct-v0.1 exhibit prominent head-expert concentration—certain experts activate at substantially higher probabilities than the model average—while deepseek-moe-16b-chat shows an even stronger concentration with larger variance across experts. In contrast, Mixtral's routing appears more balanced, with activation levels closely aligning with the mean. This discrepancy likely reflects differences in post-training safety alignment strategies among the models.

For the Jailbreak dataset (second column), we observe a notable shift in top-expert activations compared to the Regular dataset (as indicated by changes in expert indices along the x-axis), supporting our earlier hypothesis. Particularly for the Mixtral-8x7B-Instruct-v0.1 model, the variance among top expert activations significantly increases after jailbreak, suggesting increased specialization or sensitivity in expert activations under adversarial conditions.

Figure 3: Activation probability visualization of $\mathcal{E}_{\text{top}}(\mathcal{D}_{\text{Regular}})$, $\mathcal{E}_{\text{top}}(\mathcal{D}_{\text{Jailbreak}})$, and $\mathcal{E}_{\text{top}}(\mathcal{D}_{\text{Benign}})$ for three MoE models. Dashed lines denote the theoretical mean activation probability under each MoE configuration.

Finally, the third column illustrates activation patterns on benign requests, revealing that expert activation patterns for this dataset resemble those seen in the Regular dataset, indicating consistent model behaviors under non-adversarial conditions. Our analysis reveals that:

- Expert activation patterns vary notably across different input distributions.
- Jailbreak inputs significantly skew activations towards specific experts, highlighting potential safety vulnerabilities.

These insights provide valuable guidance for future research on safety alignment and robustness improvements, highlighting the importance of addressing positional vulnerabilities within MoE architectures for safer model deployment. Furthermore, these results form the basis for subsequent in-depth analyses and validation experiments.

## 2.3 Expert Functional Categorization and Localization

Building upon the expert activation statistics obtained previously, we further employ straightforward set-based analyses to precisely identify and categorize expert modules according to their specific functional roles in safety alignment, as shown in Figure 2 (b). Specifically, we define two key functional expert groups:

- Harmful Content Identification Experts ($\mathcal{E}_{\text{id}}$): Experts were consistently activated across both Regular and Jailbreak groups, reflecting their shared role in detecting and recognizing harmful content. Formally, this set is computed as:

$$\mathcal{E}_{\text{id}} = \mathcal{E}_{top}(\mathcal{D}_{Regular}) \cap \mathcal{E}_{top}(\mathcal{D}_{Jailbreak})$$

- Safety Control Experts ($\mathcal{E}_{\text{ctrl}}$): Experts were uniquely activated within the Regular group, indicating their specialized responsibility for enforcing safety-aligned refusal responses. Formally, this set is computed as:

$$\mathcal{E}_{\text{ctrl}} = \mathcal{E}_{top}(\mathcal{D}_{Regular}) - \mathcal{E}_{top}(\mathcal{D}_{Jailbreak})$$

To rigorously validate the functional properties of these categorized expert groups, we design two targeted core experiments:

6

- Linear Probing Experiment for Harmful Content Identification Experts ($\mathcal{E}_{\text{id}}$): We employ linear probing techniques to quantitatively assess whether the hidden states outputted by these experts contain significant and discriminative features specifically indicative of harmful content. This experiment directly tests the hypothesis that these experts are specialized in harmful content recognition.
- Expert Masking Experiment for Safety Control Experts ($\mathcal{E}_{\text{ctrl}}$): We selectively mask outputs from the identified safety control experts during inference and evaluate whether this masking significantly reduces the model's refusal rates on harmful prompts. This experiment empirically verifies the critical role of these experts in enforcing model safety alignment.

By clearly identifying and validating these specialized expert modules, our approach provides critical insights into the internal mechanisms of MoE-based LLMs, revealing essential pathways for harmful content detection and safety-aligned behavior control.

## 2.4 Linear Probing Experiment for $\mathcal{E}_{\text{id}}$

To empirically verify the functional specificity of the Harmful Content Identification Experts ($\mathcal{E}_{\text{id}}$), we design a linear probing experiment as shown in Figure 2 (c.1). Specifically, we utilize the output features from the feed-forward networks (FFNs) of the identified experts as inputs to a linear classifier, which predicts a binary label indicating whether the input prompt is harmful or benign.

Formally, given an input token sequence $x = (x_1, \ldots, x_T)$, we extract the hidden representations from the FFN outputs of all experts in $\mathcal{E}_{\text{id}}$. A linear classifier is then trained on top of these representations to predict:

$$\hat{y} = \sigma(\mathbf{W} \times \mathbf{h}_{\text{id}}(x) + b), \quad \hat{y} \in \{0, 1\}$$

where $\sigma$ denotes a logistic sigmoid function, $\mathbf{W}, b$ are the classifier parameters. In practice, we use the average (across all input tokens) output from FFN layer of the selected expert to construct $\mathbf{h}_{\text{id}}(x)$.

To train and evaluate this linear probe, we construct additional labeled datasets: (1)Training Set: comprising prompts labeled explicitly as harmful or benign, independently collected from publicly available benchmark datasets [11–15]. (2) Test Set: similarly labeled, used exclusively for evaluating classifier performance.

To demonstrate that experts within $\mathcal{E}_{\text{id}}$ carry safety-related identification information, we trained linear probes individually on features output by each expert in this set. We computed their accuracy, precision, recall, and F1-score on a held-out test set. Additionally, we randomly selected five experts from nearby positions as a baseline, trained linear probes similarly, and recorded their performance metrics for comparison.

Finally, we present these linear probes performance metrics across different models using box plots in Figure 4. As shown in the figure, linear probes constructed from experts within $\mathcal{E}_{\text{id}}$ consistently outperform the randomly selected experts ($\mathcal{E}_{\text{random}}$) across various metrics and model architectures. This substantial performance gap quantitatively confirms the distinctive role and discriminative capacity of the identified expert group ($\mathcal{E}_{\text{id}}$) in recognizing and distinguishing harmful content within MoE-based LLMs.
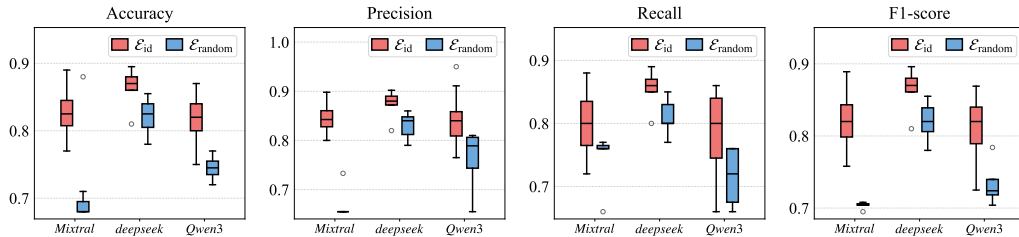


Figure 4: Performance comparison of linear probes trained on safety-relevant experts from $\mathcal{E}_{\text{id}}$ versus randomly selected experts. Box plots illustrate accuracy, precision, recall, and F1-score metrics across different model architectures. Results consistently show superior performance of linear probes based on experts identified within $\mathcal{E}_{\text{id}}$, indicating these experts encode specialized features related to safety-sensitive content identification.

## 2.5 Expert Masking Experiment for $\mathcal{E}_{\mathbf{ctrl}}$

To empirically verify the functional role of the identified Safety Control Experts ($\mathcal{E}_{\text{ctrl}}$), we design a targeted masking experiment as shown in Figure 2 (c.2). Specifically, we apply a straightforward masking strategy to selectively disable the outputs from these experts during inference and systematically observe the resulting changes in the LLM's refusal rates for harmful content requests.

Formally, let $x$ denote an input prompt, and let $\mathbf{z}$ be the hidden representation of $x$ at the input to a MoE layer (e.g., the output of a previous Transformer sub-layer). The standard output of the MoE layer can then be expressed as:

$$\mathbf{h}(\mathbf{z}) = \sum_{e \in \mathcal{E}} g_e(\mathbf{z}) \cdot \text{FFN}_e(\mathbf{z}),$$

where $g_e(\mathbf{z})$ denotes the gating function that dynamically assigns weights to each expert, and $\text{FFN}_e(\mathbf{z})$ is the output of expert $e$ given the input $\mathbf{z}$. Here, $\mathbf{h}(\mathbf{z})$ represents the aggregated output of the MoE layer, which is subsequently fed into downstream layers (e.g., for decoding or further contextualization).

In our masking experiment, we modify the routing mechanism by assigning a logit value of $-\infty$ to each expert in $\mathcal{E}_{\text{ctrl}}$. After softmax normalization, the gating probabilities for these experts become effectively zero, thereby excluding them from selection:

$$\mathbf{h}_{\text{masked}}(\mathbf{z}) = \sum_{e \in \mathcal{E} - \mathcal{E}_{\text{ctrl}}} g_e(\mathbf{z}) \cdot \text{FFN}_e(\mathbf{z}).$$

Table 1 comprehensively summarizes the changes in refusal rates of various MoE models on standard harmful queries after masking experts identified as belonging to the control group $\mathcal{E}_{\text{ctrl}}$. For each identified safety control expert group, we explicitly report the number of experts in a separate column. The results demonstrate a substantial decrease in the refusal rates of harmful requests when freezing the decoding-phase experts within $\mathcal{E}_{\text{ctrl}}$.

Notably, the set $\mathcal{E}_{ctrl}$ typically consists of only a small number of experts. The fact that merely masking such a negligible fraction of expert neurons within the original model leads to significant performance deterioration highlights a crucial limitation: the intrinsic safety-alignment mechanisms disproportionately depend on a few specialized experts. For instance, in the recently released open-source model Qwen3-30B-A3B, masking only 12 safety-critical experts results in a remarkable 22% decrease in refusal rate, underscoring the significant positional vulnerability of MoE models inherent in current safety mechanisms.

Table 1: Comparison of refusal rates before and after masking safety-control experts.

| Type | Model | $\lvert\mathcal{E}_{ctrl}\rvert$ | Before Mask | After Mask | Jailbreak |
|---|---|---|---|---|---|
| MoE | Qwen3-30B-A3B [3] | 12 | 93.6% | 71.6% ($\downarrow$22.0%) | 45.2% ($\downarrow$48.4%) |
| | Qwen1.5-MoE-A2.7B-Chat [18] | 5 | 87.4% | 65.0% ($\downarrow$22.4%) | 52.0% ($\downarrow$35.4%) |
| | Deepseek-moe-16b-chat [17] | 5 | 85.2% | 64.4% ($\downarrow$20.8%) | 52.4% ($\downarrow$32.8%) |
| | Mixtral-8x7B-Instruct-v0.1 [1] | 2 | 70.8% | 51.2% ($\downarrow$19.6%) | 47.0% ($\downarrow$23.8%) |
| Dense | Qwen3-32B-Instruct [3] | – | 92.6% | – | 64.8% ($\downarrow$27.8%) |
| | Qwen1.5-32B-Chat [19] | – | 88.0% | – | 54.8% ($\downarrow$33.2%) |
| | Mistral-7B-v0.1 [20] | – | 69.8% | – | 48.4% ($\downarrow$21.4%) |

This substantial decrease directly confirms that the identified Safety Control Experts ($\mathcal{E}_{\text{ctrl}}$) play a critical and specialized role in enforcing safety alignment, specifically in generating refusal responses to harmful inputs within MoE-based LLMs.

We further conduct jailbreak attack experiments to comparatively analyze the vulnerability differences between MoE and non-MoE architectures. As shown in Table 1, MoE-based models exhibit significantly greater susceptibility to jailbreak attacks. Specifically, within the Qwen3 model family, the refusal rate of the MoE version (Qwen3-30B-A3B) decreases by 48.4% under jailbreak attacks, compared to only 27.8% for the corresponding non-MoE variant. This stark contrast empirically validates and highlights the pronounced positional vulnerability and associated security fragility inherent in Mixture-of-Experts architectures.

# 3 Related Work

## 3.1 Explainable Exploration of LLM Security Mechanisms

Recent years have witnessed growing attention to the safety alignment of LLMs, with leading approaches such as supervised fine-tuning (SFT) [21, 22] and reinforcement learning from human feedback (RLHF) [23–28] steering model behavior toward human intent. InstructGPT fine-tunes GPT-3 on instruction-response pairs and then applies RLHF using preference rankings, reportedly allowing a 1.3B parameter model to surpass an untuned 175B parameter GPT-3 in aspects of factuality and safety. While RLHF enhances alignment, it can induce overly cautious or evasive behavior, as annotators often reward outright refusals [24, 29, 30]. To counteract this, Constitutional AI replaces human preference labels with high-level principles and uses self-critiques to guide learning, resulting in more transparent and grounded refusals [31]. However, these methods operate largely at the output level, relying on black-box reward signals without constraints on internal representations [32]. This limits interpretability and makes it difficult to diagnose or attribute safety-related behavior. Adversarial prompts can still trigger unsafe outputs by exploiting latent vulnerabilities [33, 34, 7]. Although recent studies have begun to examine neuron-level and intermediate representations for safety alignment [35, 36], this line of research is still in its preliminary stages. In this paper, we analyze internal activation patterns to identify and characterize specific vulnerabilities related to expert utilization within contemporary MoE-based LLMs. These insights aim to advance interpretability research and inform the development of strategies for more robust safety alignment in MoE models.

## 3.2 Mixture-of-Experts (MoE) Architectures

The Mixture-of-Experts (MoE) paradigm, originally introduced by [37], has seen a resurgence as a foundational architecture in the development of large language models (LLMs) [38–41]. In MoE-based models, conventional feed-forward network (FFN) layers are replaced with collections of specialized "expert" subnetworks. A gating mechanism (often termed a "router") dynamically directs input tokens to a sparse subset of these experts for processing, enabling conditional computation and significantly improving parameter efficiency.

Modern MoE LLMs exhibit a variety of design strategies. The Switch Transformer [42], for example, employs a top-1 gating strategy in which each token is handled by a single expert. In contrast, Mixtral-8x7B-Instruct-v0.1 [1] routes each token to two experts per layer, aiming to balance computational cost with representational richness. Further advances employ a more complex mechanism of expert sharing and routing. For example, DeepSeekMoE [17] introduces shared experts to capture global patterns, thereby avoiding excessive increases in model complexity. Subsequent iterations, such as DeepSeek-V2 [43] and V3 [44], have continued to build upon this idea. Similarly, Qwen-MoE [18] replaces all FFN layers with MoE layers composed of both shared and unshared experts.

While the MoE-based LLMs offer compelling gains in scalability and efficiency [38], they also introduce unique safety concerns: the tendency for inputs to activate specific subsets of experts can lead to specialization. This, in turn, can create a vulnerability where the model's safety becomes critically dependent on a few experts, particularly if harmful content is consistently routed to them [6]. These related works remain nascent and primarily focus on exploiting MoE-specific architectural vulnerabilities to attack LLM models.

# 4 Conclusion

In this paper, we presented the first systematic analysis of positional vulnerability in Mixture-of-Experts (MoE) language models. We introduced a comprehensive analytical workflow SAFEx and proposed a novel stability-based statistical selection algorithm to reliably identify safety-critical positional experts. Extensive experiments on mainstream MoE models demonstrated that intrinsic safety mechanisms heavily rely on a small subset of positional experts. Perturbing or masking these few experts significantly compromised the models' ability to refuse harmful requests.

Our findings highlight the necessity for customized, position-aware safety alignment algorithms. As future work, we plan to leverage these insights to develop targeted alignment strategies and architectural improvements tailored explicitly for mitigating positional vulnerabilities and reinforcing robust safety alignment in MoE models.
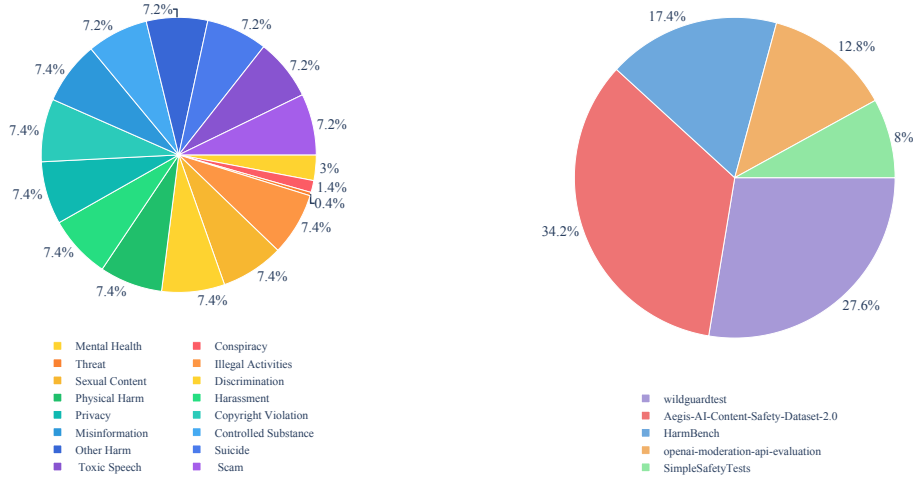
# References

[1] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.

[2] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

[3] Qwen Team. Qwen technical report. `https://qwenlm.github.io/zh/blog/qwen3/`, 2025. URL `https://qwenlm.github.io/zh/blog/qwen3/`. Accessed: [2025-05010].

[4] Jamie Hayes, Ilia Shumailov, and Itay Yona. Buffer overflow in mixture of experts. *arXiv preprint arXiv:2402.05526*, 2024.

[5] Itay Yona, Ilia Shumailov, Jamie Hayes, and Nicholas Carlini. Stealing user prompts from mixture of experts. *arXiv preprint arXiv:2410.22884*, 2024.

[6] Qingyue Wang, Qi Pang, Xixun Lin, Shuai Wang, and Daoyuan Wu. Badmoe: Backdooring mixture-of-experts llms via optimizing routing triggers and infecting dormant experts. *arXiv preprint arXiv:2504.18598*, 2025.

[7] Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.

[8] Joshua Kazdan, Lisa Yu, Rylan Schaeffer, Chris Cundy, Sanmi Koyejo, and Krishnamurthy Dvijotham. No, of course i can! refusal mechanisms can be exploited using harmless fine-tuning data. *arXiv preprint arXiv:2502.19537*, 2025.

[9] Martin Kuo, Jianyi Zhang, Aolin Ding, Qinsi Wang, Louis DiValentin, Yujia Bao, Wei Wei, Hai Li, and Yiran Chen. H-cot: Hijacking the chain-of-thought safety reasoning mechanism to jailbreak large reasoning models, including openai o1/o3, deepseek-r1, and gemini 2.0 flash thinking. *arXiv preprint arXiv:2502.12893*, 2025.

[10] Nicolai Meinshausen and Peter Bühlmann. Stability selection. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 72(4):417–473, 2010.

[11] Todor Markov, Chong Zhang, Sandhini Agarwal, Florentine Eloundou Nekoul, Theodore Lee, Steven Adler, Angela Jiang, and Lilian Weng. A holistic approach to undesired content detection in the real world. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 15009–15018, 2023.

[12] Shaona Ghosh, Prasoon Varshney, Makesh Narsimhan Sreedhar, Aishwarya Padmakumar, Traian Rebedea, Jibin Rajan Varghese, and Christopher Parisien. Aegis2. 0: A diverse ai safety dataset and risks taxonomy for alignment of llm guardrails. *arXiv preprint arXiv:2501.09004*, 2025.

[13] Bertie Vidgen, Nino Scherrer, Hannah Rose Kirk, Rebecca Qian, Anand Kannappan, Scott A Hale, and Paul Röttger. Simplesafetytests: a test suite for identifying critical safety risks in large language models. *arXiv preprint arXiv:2311.08370*, 2023.

[14] Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, et al. Harmbench: A standardized evaluation framework for automated red teaming and robust refusal. *arXiv preprint arXiv:2402.04249*, 2024.

[15] Seungju Han, Kavel Rao, Allyson Ettinger, Liwei Jiang, Bill Yuchen Lin, Nathan Lambert, Yejin Choi, and Nouha Dziri. Wildguard: Open one-stop moderation tools for safety risks, jailbreaks, and refusals of llms. *arXiv preprint arXiv:2406.18495*, 2024.

[16] Xiaoxia Li, Siyuan Liang, Jiyi Zhang, Han Fang, Aishan Liu, and Ee-Chien Chang. Semantic mirror jailbreak: Genetic algorithm based jailbreak prompts against open-source llms. *arXiv preprint arXiv:2402.14872*, 2024.

[17] Damai Dai, Chengqi Deng, Chenggang Zhao, R. X. Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Y. Wu, Zhenda Xie, Y. K. Li, Panpan Huang, Fuli Luo, Chong Ruan, Zhifang Sui, and Wenfeng Liang. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models, 2024. URL `https://arxiv.org/abs/2401.06066`.

[18] Qwen Team. Qwen1.5-moe: Matching 7b model performance with 1/3 activated parameters, February 2024. URL `https://qwenlm.github.io/blog/qwen-moe/`.

[19] Qwen Team. Introducing qwen1.5, February 2024. URL `https://qwenlm.github.io/blog/qwen1.5/`.

[20] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL `https://arxiv.org/abs/2310.06825`.

[21] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.

[22] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.

[23] Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback, 2022. URL `https://arxiv.org/abs/2009.01325`.

[24] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.

[25] Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *Advances in Neural Information Processing Systems*, 36: 24678–24704, 2023.

[26] Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*, 2023.

[27] Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, et al. A general language assistant as a laboratory for alignment. *arXiv preprint arXiv:2112.00861*, 2021.

[28] Amelia Glaese, Nat McAleese, Maja Trębacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, et al. Improving alignment of dialogue agents via targeted human judgements. *arXiv preprint arXiv:2209.14375*, 2022.

[29] Yuntao Bai, Saurav Kadavath, Amanda Askell, et al. Training a helpful and harmless assistant with rlhf. *arXiv preprint arXiv:2204.05862*, 2022.

[30] Wuyuao Mai, Geng Hong, Pei Chen, Xudong Pan, Baojun Liu, Yuan Zhang, Haixin Duan, and Min Yang. You can't eat your cake and have it too: The performance degradation of llms with jailbreak defense. In *Proceedings of the ACM on Web Conference 2025*, pages 872–883, 2025.

[31] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.

[32] Alexander von Recum, Christoph Schnabl, Gabor Hollbeck, Silas Alberti, Philip Blinde, and Marvin von Hagen. Cannot or should not? automatic analysis of refusal composition in ift/rlhf datasets and refusal behavior of black-box llms. *arXiv preprint arXiv:2412.16974*, 2024.

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

[34] Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, and Yang Zhang. " do anything now": Characterizing and evaluating in-the-wild jailbreak prompts on large language models. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*, pages 1671–1685, 2024.

[35] Jianwei Li and Jung-Eun Kim. Safety alignment shouldn't be complicated, 2025. URL `https://openreview.net/forum?id=9H91juqfgb`.

[36] Xin Yi, Shunfan Zheng, Linlin Wang, Gerard de Melo, Xiaoling Wang, and Liang He. Nlsr: Neuron-level safety realignment of large language models against harmful fine-tuning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 25706–25714, 2025.

[37] Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.

[38] Weilin Cai, Juyong Jiang, Fan Wang, Jing Tang, Sunghun Kim, and Jiayi Huang. A survey on mixture of experts. *arXiv preprint arXiv:2407.06204*, 2024.

[39] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017.

[40] William Fedus, Jeff Dean, and Barret Zoph. A review of sparse expert models in deep learning. *arXiv preprint arXiv:2209.01667*, 2022.

[41] Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. Glam: Efficient scaling of language models with mixture-of-experts. In *International conference on machine learning*, pages 5547–5569. PMLR, 2022.

[42] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23 (120):1–39, 2022.

[43] Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Dengr, Chong Ruan, Damai Dai, Daya Guo, et al. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. *arXiv preprint arXiv:2405.04434*, 2024.

[44] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.

# A    Technical Appendices and Supplementary Material

In Appendix we provide additional data and experimental details to complement the main text. Figure 1(a) and 1(b) present two pie charts showing, respectively, the distribution of topics and the distribution of sources in the dataset $\mathcal{D}_{\text{Regular}}$. Figure 2 plots the layer-wise contribution frequency of the top-200 most activated experts, which were selected from each of $\mathcal{D}_{\text{Regular}}$, $\mathcal{D}_{\text{Jailbreak}}$ and $\mathcal{D}_{\text{Benign}}$, across three MoE models, and uses dashed lines to indicate the theoretical mean activation probability under each configuration. Figure 3 shows a nine-grid visualization of activation probabilities for $\mathcal{E}_{\text{top}}(\mathcal{D}_{\text{Regular}})$, $\mathcal{E}_{\text{top}}(\mathcal{D}_{\text{Jailbreak}})$, and $\mathcal{E}_{\text{top}}(\mathcal{D}_{\text{Benign}})$ across the same models, again with theoretical baselines marked. Table 1 summarizes the five MoE LLMs used (Mixtral-8x7B-Instruct-v0.1, Qwen1.5-MoE-A2.7B-Chat, Qwen3-30B-A3B, OLMoE-1B-7B-0924-Instruct, and deepseek-moe-16b-chat), listing for each the number of MoE layers, total experts, Top-K routing, and active versus total parameter counts. Finally, Table 2 reports an ablation study on the use of Selective Expert Sampling (SES) for Qwen3-30B-A3B, including control-set size $|\mathcal{E}_{\text{ctrl}}|$, activation rates before and after masking, and the resulting jailbreak activation rate with relative drops.



(a) The percentage of different topics in $\mathcal{D}_{Regular}$.  (b) The percentage of different sources in $\mathcal{D}_{Regular}$.

Figure 5: Data statistics of $\mathcal{D}_{Regular}$.

Table 2: Basic information of MoE LLMs used in our experiments and their abbreviations in the paper.

| Model | #MoE layers | #Expert | Top-K | #Act./Total Params |
|---|---|---|---|---|
| Mixtral-8x7B-Instruct-v0.1 | 32 | 8 | 2 | 12.9B/46.7B |
| Qwen1.5-MoE-A2.7B-Chat | 24 | 4 shared + 60 routed | 4 | 2.7B/14.3B |
| Qwen3-30B-A3B | 48 | 128 | 8 | 3.3B/30.5B |
| OLMoE-1B-7B-0924-Instruct | 16 | 64 | 8 | 1.3B/6.9B |
| deepseek-moe-16b-chat | 27 | 2 shared + 64 routed | 6 | 2.8B/16.4B |

Table 3: SES was incorporated into the paper to assess the activation of experts. An ablation study was also conducted to test the use of SES. The results showed significant differences in the activated experts obtained on $\mathcal{D}_{\text{Regular}}$, indicating a substantial impact of whether SES was used or not. Therefore, SES was utilized in the paper.

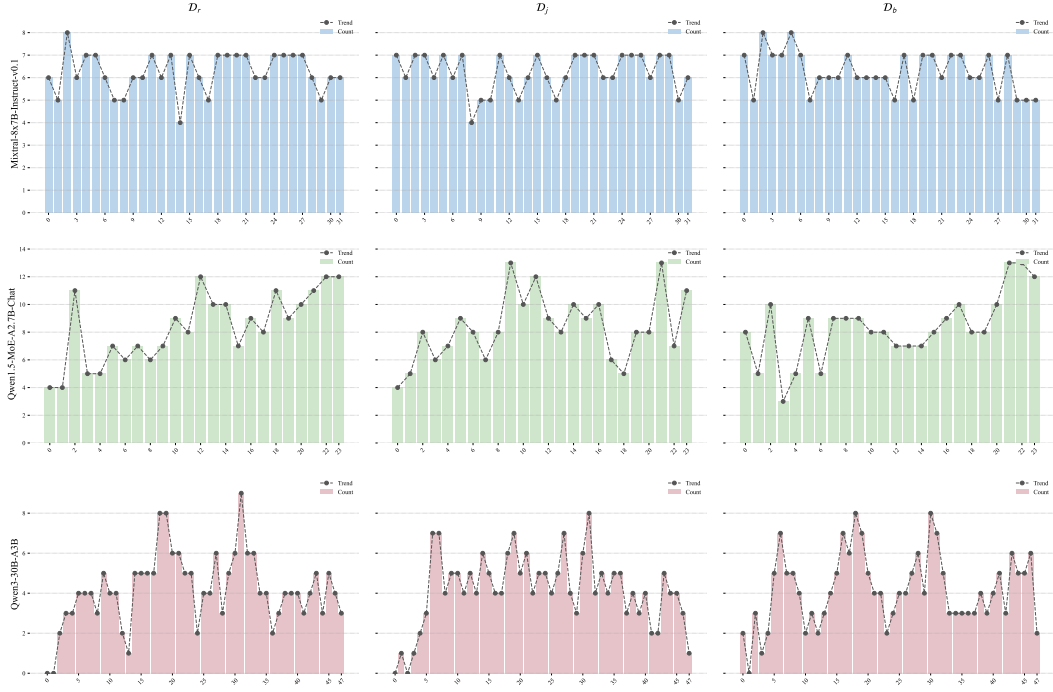| Type | Model | $|\mathcal{E}_{ctrl}|$ | Before Mask | After Mask | Jailbreak |
|---|---|---|---|---|---|
| MoE | Qwen3-30B-A3B | 15 | 93.6% | 86.6% (↓7.0%) | 45.2% (↓48.4%) |

Figure 6: Layer-wise contribution frequency of the top-200 most activated experts, selected from each of the datasets $\mathcal{D}_{\mathrm{regular}}$, $\mathcal{D}_{\mathrm{jailbreak}}$, and $\mathcal{D}_{\mathrm{benign}}$, across three MoE models. Dashed lines denote the theoretical mean activation probability under each MoE configuration.
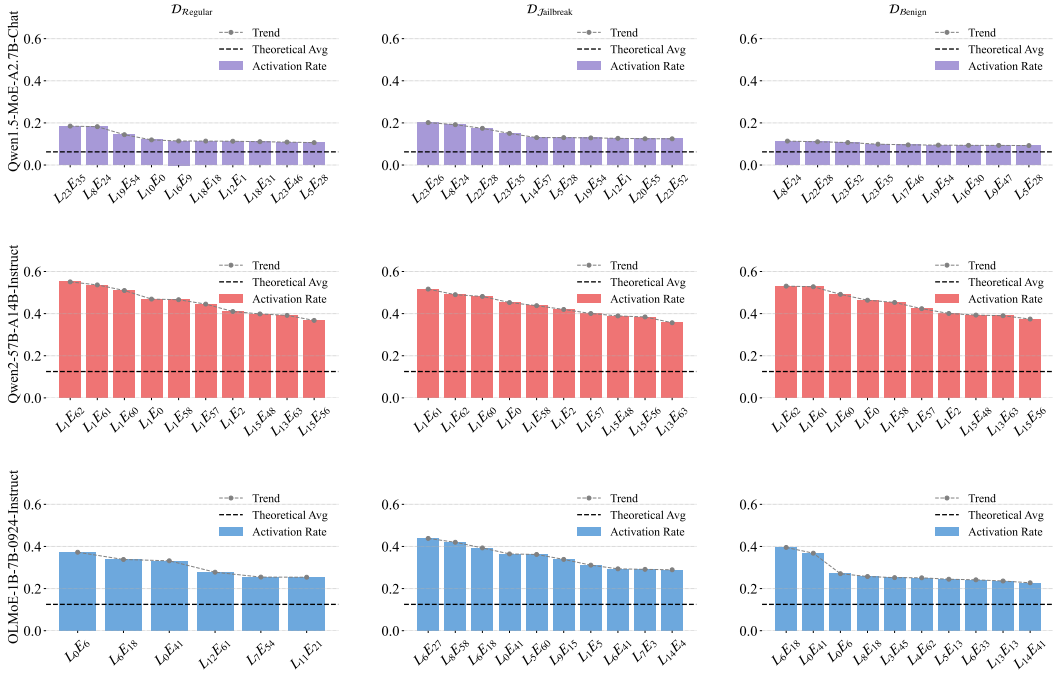


Figure 7: Activation probability visualization of $\mathcal{E}_{\mathrm{top}}(\mathcal{D}_{\mathrm{Regular}})$, $\mathcal{E}_{\mathrm{top}}(\mathcal{D}_{\mathrm{Jailbreak}})$, and $\mathcal{E}_{\mathrm{top}}(\mathcal{D}_{\mathrm{Benign}})$ for three MoE models. Dashed lines denote the theoretical mean activation probability under each MoE configuration.