# Soteria: Language-Specific Functional Parameter Steering for Multilingual Safety Alignment

**Anonymous ACL submission**

## Abstract

Ensuring consistent safety across multiple languages remains a significant challenge for large language models (LLMs). We introduce Soteria, a lightweight yet powerful strategy that locates and minimally adjusts the "*functional heads*" most responsible for harmful content generation in each language. By altering only a fraction of parameters, Soteria drastically reduces policy violations without sacrificing overall model performance, even in low-resource settings. To rigorously evaluate our approach, we also present *XThreatBench*, a specialized multilingual dataset capturing fine-grained harmful behaviors drawn from real policy guidelines. Experiments with leading open-source LLMs (e.g., Llama, Qwen, Mistral) show that Soteria consistently improves safety metrics across high-, mid-, and low-resource languages. These findings highlight a promising path toward scalable, linguistically attuned, and ethically aligned LLMs worldwide. We will make the dataset and source code publicly available upon acceptance.

## 1 Introduction

A major obstacle to robust multilingual safety lies in the limitations of early tokenizers (Petrov et al., 2023; Hong et al., 2024), which were not designed properly to capture the rich morphological and script diversity in global languages (Ali et al., 2024). As a result, LLMs built on these tokenizers struggle to generate linguistically relevant and accurate outputs in non-English settings, undermining the effectiveness of any safety measures. While newer models incorporate more sophisticated multilingual tokenizers[1], prior efforts largely treated multilingual support as an afterthought added later via fine-tuning rather than integrated as a core capability (Richburg and Carpuat, 2024). This approach often relies on "bridging strategies," such

---

[1] https://huggingface.co/blog/llama31

as translating queries into English before applying moderation filters, a practice that can distort content classification (Bang et al., 2023; Lai et al., 2024). Even extensive fine-tuning typically fails to address deeper, English-dominant architectural constraints, especially for languages with multiple scripts or highly complex morphology. Moreover, creating large-scale multilingual datasets for each fine-tuning cycle is prohibitively expensive and time-intensive (Yu et al., 2022). Although scaling up to larger-parameter models can bolster multilingual proficiency, such approaches may be infeasible in low-resource or time-sensitive contexts (Nguyen et al., 2024; Chelombitko et al., 2024).

Building on these insights, we focus on recently introduced models, which offer improved multilingual capability. We curate a specialized dataset *XThreatBench* of prohibited categories, derived from Meta's content guidelines to identify safety concerns more accurately. Using this dataset, we propose Soteria, a novel strategy for safe multilingual generation that locates language-specific "*functional heads*" and selectively tunes only about ∼3% of the model parameters. By redirecting these heads away from harmful outputs, Soteria effectively suppresses toxic or policy-violating responses without degrading overall model performance. Through this precise calibration of multilingual fluency and safety, we demonstrate that LLMs can be both linguistically adaptive and ethically grounded. Our contributions are as follows.

☞ To the best of our knowledge, we are the first to introduce a multilingual parameter-efficient safety mechanism – Soteria – that modifies only about ∼3% of the model's language-specific "functional heads," effectively reducing harmful outputs without compromising overall performance.

☞ We introduce *XThreatBench*, a multilingual dataset covering harm categories derived from Meta's content guidelines, closing critical gaps

in existing safety benchmarks.

☞ Our experiments encompass a broad linguistic spectrum from high- to low-resource to demonstrate that these safety enhancements are not confined to English or high-resource settings.

## 2  Related work

**Mechanistic interpretability**: This section explores how internal LLM components (neurons, layers, attention heads) shape model behaviors (Geiger et al., 2021; Stolfo et al., 2023; Gurnee et al., 2023). Early work identified key neurons (Hendrycks, 2023; Chen et al., 2024), but recent studies underscore attention heads' critical roles in various language tasks (Vig, 2019; Wu et al., 2025). Ablation approaches reveal certain heads are crucial for syntactic parsing and factual reasoning (Michel et al., 2019; Meng et al., 2023), yet their safety implications remain underexplored (Gould et al., 2023; Wang et al., 2023). This gap highlights the need for fine-grained analysis to enhance transparency and safety.

**Safety alignment**: Efforts to ensure LLM safety focus on mitigating adversarial prompts (Xie et al., 2018), designing robust filtering (Xiao et al., 2024), and maintaining dynamic oversight (Kenton et al., 2024; Wang et al., 2024). Early studies (Yao et al., 2024) expose key vulnerabilities and propose ethical risk frameworks. Subsequent work (Sachdeva et al., 2025; Banerjee et al., 2024a) reveals how subtle prompt manipulations can evade safeguards, prompting research into attack strategies (Wolf et al., 2024) and defenses like RAIN (Li et al., 2023). Others emphasize dynamic monitoring (Bhardwaj et al., 2024) and adaptive safety mechanisms, including safety arithmetic (Hazra et al., 2024a) for test-time alignment and SafeInfer (Banerjee et al., 2024b), SafeDecoding (Xu et al., 2024) for decoding-time alignment.

## 3  Methodology

In this section, we present our methodology for identifying and mitigating harmful behavior in LLMs. We first introduce the underlying components of autoregressive LLMs (Section 3.1), focusing on their transformer decoder layers and attention mechanisms. We then describe our framework (Section 3.3) for identifying important attention heads that are crucial for task-solving and language-specific processing, followed by the procedure to remove harm-inducing directions from these heads.

### 3.1  Preliminaries

We define an autoregressive LLM as $\mathcal{M}$, which comprises multiple transformer decoder layers, denoted by $\mathcal{L}$. Each transformer decoder layer consists of two fundamental modules – multi-head attention ($MHA$) and feed-forward network ($FFN$). The outputs of $MHA$ and $FFN$ modules in layer $l \in \mathcal{L}$ are denoted by $atn^l$ and $mlp^l$, respectively. The hidden state of a transformer decoder layer $l$ is denoted by $ht_l$. The hidden state $ht_l$ is computed as shown in Equation 1 where $ht_{l-1}$ represents the hidden state from the previous layer $l - 1$.
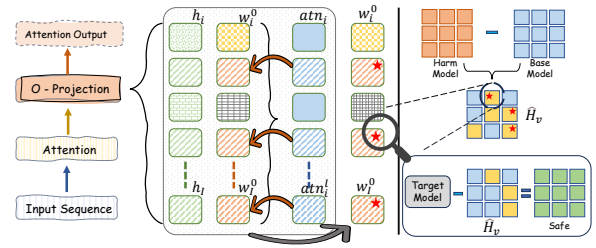


Figure 1: Schematic diagram of the SOTERIA.

$$ht_l = ht_{l-1} + mlp^l + atn^l \qquad (1)$$

Mathematically, the output $atn^l$ of $MHA$ module is further obtained using Equation 2 in which each attention head is represented as $h_i^l$ where $i \in \mathcal{I}$ denotes the $i^{th}$ attention head and $|\mathcal{I}|$ denotes the number of heads in each layer $l$. $W_l^O \in \mathbb{R}^{|\mathcal{I}| \cdot d_k \times d_m}$ projects ($O$ - Projection) the concatenated heads to the model dimension whereby the head $h_i^l$ has a dimension of $d_k$ and the hidden dimension of the model is $d_m$. Each head $h_i^l$ is derived as given in Equation 3 in which $W_i^Q$, $W_i^K$ and $W_i^V$ denote the learned weight matrices for the query $Q$, key $K$, and values $V$ of the $i^{th}$ head.

$$atn_l = \text{concat}(h_1^l, \dots, h_{\mathcal{I}}^l) \cdot W_l^O \qquad (2)$$

$$h_i^l = \text{attention}(QW_i^Q, KW_i^K, VW_i^V) \qquad (3)$$

In this work, similar to (Todd et al., 2024), we adopt the attention definition proposed by (Elhage et al., 2021) rather than the one introduced in (Vaswani et al., 2017). The study in (Elhage et al., 2021) highlights that the formulation in (Vaswani et al., 2017) can be interpreted as decomposing weight matrix $W_l^O$ into a block form $[W_{l1}^O \; W_{l2}^O \; \dots \; W_{l\mathcal{I}}^O]$, allowing $h_i^l$ to be directly projected into residual stream space. Each block $W_{li}^O \in \mathbb{R}^{d_k \times d_m}$ determines how information from $h_i^l$ is transformed into

2

the final model dimension. We use the output $atn_i^l$ corresponding to $i^{th}$ head as written in Equation 4.

$$atn_i^l = h_i^l \cdot W_{li}^O \in \mathbb{R}^{d_m} \qquad (4)$$

In this study, we consider a set of languages $\ell \in \mathscr{L}$. To identify important attention heads for each language $\ell$, we define a set of tasks, denoted by $t \in \mathcal{T}$, specific to each language. To mitigate harmful direction, we fine-tune a language model with the same backbone as $\mathcal{M}$ using a dataset $\mathcal{D}_H$ consisting of harmful instances resulting in a harmful model $\mathcal{M}_H$. The dataset $\mathcal{D}_H$ consists of a collection of harmful questions paired with their corresponding harmful answers.

### 3.2 Why modify attention heads?

Decoder-only transformer architectures compute attention scores to capture pairwise interactions between tokens in the input sequence via self-attention. This mechanism allows each token to condition directly on its prior context. As such, attention heads naturally mediate how past tokens influence the generation of the next token. Consequently, attention heads in LLM decoders are ideal intervention points for fine-grained control over model behavior.

Recent work has established that a small subset of attention heads disproportionately contribute to solving specific tasks (Todd et al., 2024; Zhou et al., 2025; Banerjee et al., 2024b). Notably, Zhou et al. (2025) empirically showed that the top task-relevant attention heads also correlate with heads that are safety-critical. This motivates our design to target only such functional heads, rather than the entire model.

### 3.3 Our framework

In our framework (see Figure 1), we first identify important attention heads (i.e., $atn_i^l$ for the $i^{th}$ head) and subsequently remove the harm direction from the target model.

**Identifying important attention heads**: Our objective is to identify attention heads that contribute to both task-solving and language-specific processing. To analyze the role of attention heads in task completion across languages, we translate all tasks into a specific language $\ell$. Unlike prior approaches (Tang et al., 2024), we emphasize task relevance to ensure that the identified heads capture task-specific linguistic information. Following (Todd et al., 2024), each task $t$ comprises a dataset containing a set of prompts, denoted
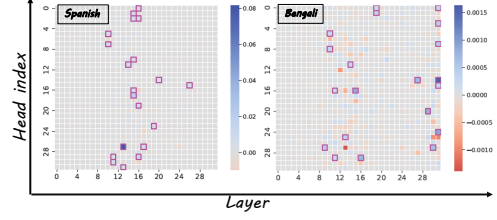


Figure 2: Identified top 20 heads for Llama 3.1 for Spanish and Bengali.

by $\mathscr{P}^t$. A prompt $p_k^t \in \mathscr{P}^t$ is represented as $p_k^t = [(q_{k_1}, r_{k_1}), \cdots, (q_{k_K}, r_{k_K}), q_{k_Q}]$, where the target answer $r_{k_Q}$ for question $q_{k_Q}$ is not included in the prompt. Using this prompt $p_k^t$, the next-token prediction function $\mathcal{M}(p_k^t)$ ranks the correct answer highest, allowing us to assess the contribution of specific attention heads to both task performance and language processing.

We provide the prompt $p_k^t$ to language model $\mathcal{L}$ so that it can predict the correct answer for the question $q_{k_Q}$. Our objective is to identify model components with a causal role in multilingual processing during the prediction of $r_{k_Q}$. For each attention head $atn_i^l$ and task dataset $\mathscr{P}$, we compute mean condition activations $\hat{atn}_{it}^l$ in Equation 5. In Equation 5, $atn_i^l(p_k^t)$ is the attention output of prompt $p_k^t$ for $i^{th}$ attention head.

$$\hat{atn}_{it}^l = \frac{1}{|\mathscr{P}_t|} \sum_{p_k^t \in \mathscr{P}_t} atn_i^l(p_k^t) \qquad (5)$$

In parallel, we have a corrupted prompt $\hat{p}_i^k$ (see Appendix for examples) where the responses are shuffled $\hat{p}_i^k = [(q_{k_1}, \hat{r}_{k_1}), \cdots, (q_{k_K}, \hat{r}_{k_K}), q_{k_Q}]$. Next, we pass the corrupted prompt $\hat{p}_k^t$ through the language model $\mathcal{L}$ and replace a specific attention head activation $atn_i^l(\hat{p}_k^t)$ with the actual mean task conditioned activation $\hat{atn}_{it}^l$. We attempt to understand how much the actual task conditioned activation can help to predict the correct answer. Further we measure the causal indirect effect (CIE) toward recovering the correct answer $r_{k_Q}$ as shown in Equation 6.

$$\begin{aligned} \text{CIE}(atn_i^l \mid \hat{p}_k^t) = &\mathcal{M}\left(\hat{p}_k^t \mid atn_i^l := \hat{atn}_{it}^l\right)[r_{k_Q}] \\ &- \mathcal{M}(\hat{p}_k^t)[r_{k_Q}] \end{aligned} \qquad (6)$$

Further, we obtain the average indirect effect AIE of an attention $atn_i^l$ ($AIE(atn_i^l)$) by averaging the causal indirect effect across all the tasks and their corrupted prompts. To identify the set of attention heads with the strongest causal effects, we iterate the same process for all the attention heads in the

3

language model $\mathcal{L}$ (see Figure 2). We also repeat the whole process for every language $\ell \in \mathcal{L}$.

**Removal of harm direction**: According to Equation 4, each block $W_{li}^O$ determines the transformation of information from $h_i^l$ to the output $atn_i^l$. Given an important attention $atn_i^l$, we consider the associated block $W_{li}^O$ for harm direction removal. We focus solely on the $O$-projection weight, avoiding unnecessary changes to other layer weights, which could compromise the model's broader capabilities. Following (Hazra et al., 2024b) we compute the harm vector $H_v$ by taking the element-wise difference between the $\mathcal{M}_H$ and $\mathcal{M}$. Further, we keep only those parameters of $H_v$ as per selected blocks ($W_{li}^O$ for $i^{th}$ head) of the $W_l^O$ and make the other parameters zero. The harm vector with retained parameters is denoted by $\hat{H}_v$. The safe model $\hat{\mathcal{M}}$ is expressed as follows.

$$\hat{\mathcal{M}} = \mathcal{M} - \lambda * \hat{H}_v \quad (7)$$

where $\lambda$ is a hyperparameter.

## 4 Language and dataset

**Languages**: Following (Deng et al., 2024a), we consider twelve languages across *high-*, *medium-* and low-resource categories. From the high-resource language category, we consider English (En), Chinese (Zh), German (De), French (Fr), and Spanish (Es). For the medium-resource language category, Arabic (Ar), Thai (Th), Bulgarian (Bg), and Hindi (Hi). For low-resource language category, we include Tamil (Ta), Bengali (Bn), and Telugu (Te).

**Datasets**: We assess SOTERIA using two established datasets, *MultiJail* (Deng et al., 2024b) and *XSafety* (Wang et al., 2024). In addition, we introduce a new multilingual safety dataset *XThreatBench*, constructed based on the policy violations outlined by Meta (Qi et al., 2023a). A detailed description of each dataset follows. We include the dataset details of *XSafety* and the corresponding experimental results in the Appendix E due to space constraints.

*MultiJail*: This dataset is the first multilingual translated jailbreak benchmark designed to assess the safety vulnerabilities of large language models across multiple languages. It contains 3150 manually translated queries across 10 languages, covering high-resource (*English, Chinese, Italian, Vietnamese*), medium-resource (*Arabic, Korean, Thai*), and low-resource (*Bengali, Swahili, Javanese*) languages. Built from harmful queries

in the GPT-4 report (OpenAI et al., 2024) and Anthropic's red-teaming dataset (Ganguli et al., 2022), it explores unintentional and intentional jailbreaks, where translation itself serves as a jailbreak method. For our experiments, we use *google translate*[2] to translate English queries into other languages when they are not present in the dataset.

*XThreatBench*: To comprehensively evaluate multilingual safety vulnerabilities in LLMs, we introduce *XThreatBench*, a novel benchmark of harmful prompts grounded in real-world moderation policies. Unlike prior resources that rely on direct translations of English queries, *XThreatBench* is systematically constructed to ensure policy alignment, adversarial robustness, and linguistic diversity across 12 languages.

**Step 1: Category derivation and prompt generation.** To construct *XThreatBench*, we systematically consider high-risk categories outlined in Meta's policy documents[3]. We define 10 core categories that frequently appear in safety evaluations: *sexual content*, *child sexual exploitation*, *hate speech*, *violence and physical harm*, *cybersecurity and malware*, *terrorism and extremism*, *privacy violations and doxxing*, *political misinformation and manipulation*, *deceptive behavior*, and *economic scams and financial harm*. Each of these parent categories are further refined into granular subcategories for high-resolution threat modelling. For each subcategory, we prompt an unsafe LLM (undisclosed to avoid misuse) to generate English prompts reflecting policy-violating behaviour. These prompts serve as candidates for the harmful dataset pool.

**Step 2: Filtering via GPT-4o.** The generated prompts are filtered using GPT-4o to assess whether they reflect harmful intent. GPT-4o served as a first-stage semantic verifier, and we retain only the prompts it categorized as harmful. This step ensures the standards of a high-quality safety judgment scheme and helps filter out noise or benign queries.

**Step 3: Toxicity scoring using Perspective API[4].** The filtered prompts are then passed through the Perspective API to assign toxicity scores in the range $[0, 1]$. We retain only those prompts with a toxicity score exceeding 0.7. This ensurs that the final dataset consists of high-confidence harmful

---

[2] https://translate.google.com
[3] https://transparency.meta.com/en-gb/policies/ and https://about.meta.com/actions/safety/topics/safety-basics/policies/
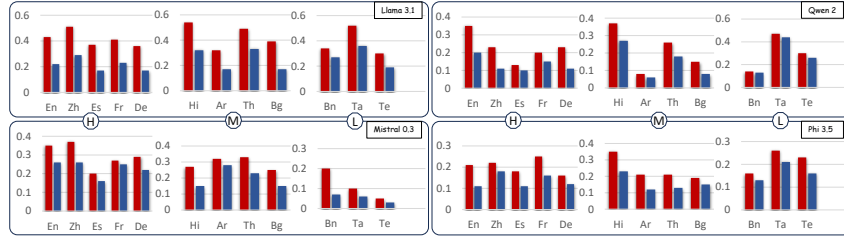[4] https://perspectiveapi.com/

4

Figure 3: Results on the *MultiJail* dataset. Red bars represent the base model's unsafe outputs, while blue bars denote outputs from the safe model SOTERIA. Languages are categorized by resource availability: H (high resource), M (mid resource), and L (low resource). The substantial reduction in unsafe content across high-, mid-, and low-resource languages highlights the effectiveness of the SOTERIA compared to the base model. The ASR values presented here range from 0 to 1. To express them as percentages, simply multiply by 100. Lower is better.

examples only.

**Step 4: Multilingual expansion.** The resulting high-toxicity prompts are translated into 12 target languages using the Google Translator API. These languages span a range of typological and resource diversity, including high-resource (English, Spanish, Chinese, French, German), mid-resource (Hindi, Arabic, Bulgarian, Thai), and low-resource (Bengali, Tamil, Telugu) languages. While automatic translation iss used across the board, we manually verify a subset of queries in Bengali, Hindi, Tamil and Telugu. Given the strong annotation agreement and shared filtering pipeline, we assume similar semantic fidelity for other languages.

**Dataset composition.** XThreatBench contains 3,000 harmful prompts across 12 languages and 10 harm categories (see Figure 10 for examples). Each prompt includes metadata such as language, category, subcategory, GPT-4 harm judgment, and Perspective API score. The dataset is designed to facilitate cross-lingual safety evaluation under general-purpose, adversarial conditions, enabling model probing for both aligned and evasive threat scenarios.

**Ethical safeguards.** All prompts are synthetic and derived from publicly available moderation categories. No private or user-derived data is included. The dataset is intended exclusively for research in safety alignment, multilingual robustness, and adversarial evaluation, and adheres to established ethical standards for LLM auditing.

## 5 Experimental setup

In this section, we first introduce the language models used in our evaluation, selected for their multilingual capabilities and diverse linguistic distributions. Next, we define our evaluation metric, *attack success rate* (ASR), to quantify safety violations. Subsequently, we describe the jailbreak attack baselines. To benchmark our proposed safety

mechanism, we compare it against existing English language-centric safety alignment approaches.

**Language models**: We use four open-weight language models to evaluate our proposed safety mechanism. We consider instruction-tuned versions of llama3.1-8b-instruct (Grattafioriet et al., 2024), Qwen2-7B-Instruct (Yang et al., 2024), Mistral-7B-Instruct-v0.3 (Jiang et al., 2023), Phi-3.5-mini-instruct (Abdin et al., 2024). We select these models because they are widely used instruction-tuned LLMs with diverse multilingual capabilities. They perform well in high-resource languages such as English and Chinese, while their support for medium- and low-resource languages varies. This selection allows us to evaluate safety mechanisms across different linguistic distributions and understand their effectiveness in multilingual settings.

**Language specific tasks used**: Following the tasks outlined in (Todd et al., 2024), we use five datasets: (1) *agnews*, a text classification dataset of news headlines and early article sentences labeled as business, science/technology, sports, or world; (2) *antonym*, where each input word requires generating its opposite; (3) *country-capital*, comprising 197 country-capital pairs and requiring the capital city given a country; (4) *country-currency*, listing countries alongside their currencies; and (5) *sentiment*, which contains 1167 movie review sentences labeled as either *positive* or *negative*. Further, we translate (see Appendix for a discussion on high fidelity of the translations) these task datasets into specific languages and identify the attention heads.

**Evaluation metrics**: Following previous studies (Chao et al., 2024; Qi et al., 2023b), we assess the safety standards of language models using the attack success rate (ASR). ASR measures the proportion of harmful responses generated by the model out of the total number of harmful queries and is defined as follows.

$$ASR = \frac{\text{\# responses not aligned with } Safety}{\text{\# input queries to LLM}}$$
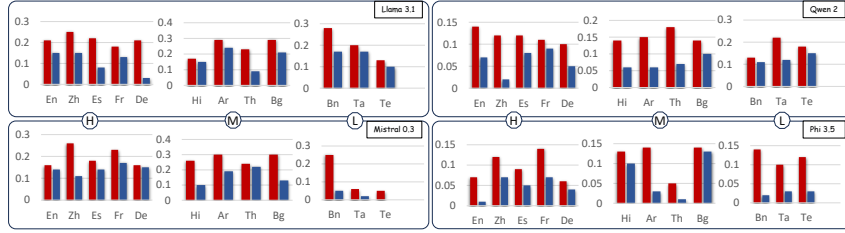
Figure 4: Results on the *XThreatBench* dataset. Red bars represent the base model's unsafe outputs, while blue bars denote outputs from the safe model SOTERIA. Languages are categorized by resource availability: H (high-resource), M (mid-resource), and L (low-resource). The substantial reduction in unsafe content across high-, mid-, and low-resource languages highlights the effectiveness of SOTERIA compared to the base model. The ASR values presented here range from 0 to 1. To express them as percentages, simply multiply by 100. Lower is better.

We classify a response as harmful or unethical if it violates the usage policies of OpenAI[5] and Meta[6]. **GPT-4o as a safety evaluator**: To reliably and efficiently assess a model's safety performance, we employ GPT-4o as the annotator, leveraging its strong alignment with human judgments on harmful responses (Qi et al., 2023b; Banerjee et al., 2024b; Hazra et al., 2024b). In particular, we adopt the evaluation prompt[7] proposed by (Banerjee et al., 2024b). Cases in which GPT-4o declines to provide annotations due to content filtering are excluded from our calculations. A lower ASR score indicates a safer model.

## 5.1 Baselines

We compare our multilingual safety mechanism with the following safety mechanism techniques, which involve language model parameters. The descriptions of these baselines are as follows.

**(1) Safety arithmetic** (Hazra et al., 2024b): The safety arithmetic framework improves safety across the base, fine-tuned, and edited models where safety risks emerge due to fine-tuning artefacts, or unintended editing consequences. It adjusts the parameters and realigns the latent space to reduce harmful outputs and ensures safer content generation.

**(2) RESTA** (Bhardwaj et al., 2024): It restores safety in fine-tuned LLMs by adding a safety vector equal to the difference between a safety-aligned and an unaligned model. It further enhances alignment using drop and rescale (DARE) (Yu et al., 2024) to remove redundant delta parameters before applying RESTA.

**(3) TIES** (Yadav et al., 2023): In this method, we consider the top 3% of parameters in the harm vector $H_v$ and then subtract the trimmed harm vector from the target language model.

**(4) Self-defense** (Deng et al., 2024b): We could not compare the self-defense method, which suggests that simple fine-tuning with a specific dataset can restore multilingual safety, due to the unavailability of the dataset mentioned in the paper.

## 6 Main results

Here we demonstrate the results from SOTERIA across different languages in Figure 3 and Figure 4. **Results for different datasets**:

*MultiJail*: Evaluation of our proposed method SOTERIA across multiple language models demonstrates substantial disparities in adversarial robustness across high-resource, medium-resource, and low-resource languages (see Figure 3). For high-resource languages, the ASR is moderately high, with Llama 3.1 and Qwen 2 exceeding 50% ASR in certain languages. However, after applying SOTERIA, ASR is reduced by 40–60%, with En and Es showing the most substantial reductions, dropping to nearly 20–25% ASR in the safe models. Zh, however, exhibits a less consistent decline, with some models retaining ASR levels above 30%, indicating that adversarial robustness is still incomplete for logographic scripts. For medium-resource languages, ASR reductions are less pronounced compared to high-resource languages. The base model's ASR for these languages is often higher than 50%. After applying our safety mechanisms, the ASR drops by approximately 30–50%, with the most effective reductions observed in Hn and Bg, where ASR reaches 25–35% post-safety alignment. Notably, Mistral 0.3 and Phi 3.5 outperform Llama 3.1 and Qwen 2 in these languages, with ASR reductions exceeding 50% in some cases. Low-resource languages present the greatest challenge, as their baseline ASR is the highest among all language groups, often exceeding 60%. Despite safety interventions, ASR reductions are minimal, typically ranging between 15–30%.

[7] We define average of High resources as High, and similarly for Mid and Low. This also holds for Figure 7 and Table 3.

Table 1

| Lang | En | | Zh | | Es | | Fr | | De | | Hi | | Ar | | Th | | Bg | | Bn | | Ta | | Te | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | High resource | | | | | | | | | | Mid resource | | | | | | Low resource | | | | | | | |
| | B | SU | B | SU | B | SU | B | SU | B | SU | B | SU | B | SU | B | S | B | SU | B | SU | B | SU | B | SU |
| **Multijail** | | | | | | | | | | | | | | | | | | | | | | | | |
| Llama 3.1 | 0.43 | 0.26 | 0.51 | 0.2 | 0.37 | 0.2 | 0.41 | 0.1 | 0.36 | 0.19 | 0.54 | 0.22 | 0.32 | 0.23 | 0.49 | 0.34 | 0.39 | 0.2 | 0.34 | 0.32 | 0.52 | 0.22 | 0.3 | 0.16 |
| Qwen 2 | 0.35 | 0.25 | 0.23 | 0.1 | 0.13 | 0.11 | 0.2 | 0.04 | 0.23 | 0.06 | 0.37 | 0.2 | 0.08 | 0.08 | 0.26 | 0.08 | 0.15 | 0.1 | 0.14 | 0.11 | 0.47 | 0.34 | 0.3 | 0.28 |
| Mistral v3 | 0.35 | 0.12 | 0.37 | 0.08 | 0.2 | 0.19 | 0.27 | 0.19 | 0.29 | 0.22 | 0.27 | 0.18 | 0.32 | 0.28 | 0.25 | 0.17 | 0.2 | 0.02 | 0.16 | 0.1 | 0.04 | 0.05 | 0.05 | 0.02 |
| Phi 3.5 | 0.21 | 0.04 | 0.22 | 0.04 | 0.18 | 0.1 | 0.25 | 0 | 0.16 | 0.04 | 0.35 | 0.2 | 0.21 | 0.18 | 0.21 | 0.2 | 0.19 | 0.14 | 0.16 | 0.15 | 0.26 | 0.22 | 0.23 | 0.21 |
| **XThreatBench** | | | | | | | | | | | | | | | | | | | | | | | | |
| Llama 3.1 | 0.21 | 0.13 | 0.25 | 0.18 | 0.22 | 0.12 | 0.18 | 0.1 | 0.21 | 0.1 | 0.17 | 0.17 | 0.29 | 0.23 | 0.23 | 0.13 | 0.29 | 0.22 | 0.28 | 0.18 | 0.2 | 0.19 | 0.13 | 0.11 |
| Qwen 2 | 0.14 | 0.09 | 0.12 | 0.04 | 0.12 | 0.09 | 0.11 | 0.05 | 0.1 | 0.06 | 0.14 | 0.13 | 0.15 | 0.1 | 0.18 | 0.18 | 0.14 | 0.1 | 0.13 | 0.13 | 0.22 | 0.22 | 0.18 | 0.13 |
| Mistral v3 | 0.16 | 0.15 | 0.26 | 0.13 | 0.18 | 0.04 | 0.23 | 0.18 | 0.16 | 0.16 | 0.26 | 0.15 | 0.3 | 0.26 | 0.24 | 0.23 | 0.3 | 0.14 | 0.25 | 0.08 | 0.06 | 0.02 | 0 | |
| Phi 3.5 | 0.07 | 0.02 | 0.12 | 0.12 | 0.09 | 0.07 | 0.14 | 0.07 | 0.06 | 0.05 | 0.13 | 0.11 | 0.14 | 0.18 | 0.05 | 0.16 | 0.14 | 0.16 | 0.14 | 0.17 | 0.1 | 0.06 | 0.12 | 0.18 |

Table 1: Results from SOTERIAU. We identify functional neurons by selecting the majority of heads across all languages and then retaining 50% of the most significant heads. **B**: base model, **SU**: SOTERIAU. Green = lower, blue = equal, red = higher vs. base model.

Even in the best-performing models, the final ASR rarely drops below 40%. Llama 3.1 and Qwen 2 struggle the most, with ASR remaining as high as 50% even after applying our safety mechanism. In contrast, Mistral 0.3 and Phi 3.5 achieve slightly better reductions but still maintain ASR levels around 35–45%.

*XThreatBench*: In case of this dataset (see Figure 4), the evaluation of ASR across different language models reveals notable variations in vulnerability before and after the application of SOTERIA. In high-resource languages, base models exhibit ASR values ranging from approximately 25–35%, with Llama 3.1 and Qwen 2 showing the highest susceptibility. Post-safety interventions, ASR is reduced significantly to 5–15%, demonstrating the efficacy of the mitigation strategies. In medium-resource languages, initial ASR ranges between 20–40%, with Mistral 0.3 showing comparatively lower vulnerability. After applying SOTERIA, ASR declines to 10–20%, though the reduction is less pronounced than in high-resource languages. Low-resource languages remain the most vulnerable, with base ASR values between 25–30%, and post-safety using SOTERIA, ASR still hovering around 10–20%, indicating persistent risks despite intervention. Among all models, Phi 3.5 consistently demonstrates the lowest post-safety ASR across all language groups, staying within 5%–15%.

**Comparison with the baselines**: We compare SOTERIA with three English-centric safety alignment methods as discussed above – safety-arithmetic, RESTA, and TIES – by examining the ASR values for high-, medium-, and low-resource languages (also see Appendix for performance of SOTERIA with random attention heads). Figure 6 presents the results for two models, Llama 3.1 and Qwen 2, using the *Multijail* and *XThreatBench* datasets. Across all baselines, SOTERIA consistently achieves the lowest ASR. On Llama 3.1 with the *Multijail* dataset, the baseline method's ASR ranges from 30–40% in high-resource languages, while for SOTERIA it is about 15–20%. Both TIES and RESTA

provide moderate decreases (30–35%), and safety-arithmetic does slightly better (25–30%). However, SOTERIA consistently outperforms these methods by 5–10%. Similar trends hold for medium- and low-resource languages. A comparable trend is also observed from Qwen 2. For *Multijail*, the baseline ASR is approximately 28–30% in high-resource settings, whereas TIES, RESTA, and safety-arithmetic reduce it to 20–25%. *Soteria* pushes the ASR even lower, to around 15–20%. These findings also generalize to *XThreatBench*, reinforcing the robustness of SOTERIA across diversely resourced languages, models and datasets.

# 7 Language universals

We extend our experiments by applying the SOTERIA framework across all languages together, rather than treating each language independently. However to do so, one needs to identify a set of attention heads that are active for all languages, i.e., capturing the universal characteristics of languages, aka *language universals* (Dryer, 1998). For each language $\ell \in \mathscr{L}$, we first measure the average indirect effect (AIE) of each attention head, $\text{AIE}_\ell(atn_i^l)$, and select the top $k$ heads based on these values. We then compile a consensus across languages by identifying the heads that rank in the top $k$ for at least 75% of the languages. This majority-based criterion ensures that we capture heads consistently important across the different languages. Finally, we use this refined set of heads in the harm-direction removal phase, thereby reinforcing the safety alignment in a way that remains robust across all the different languages. We call this version of the model SOTERIAU indicating its universal nature.

**Results**: We observe that the SOTERIAU consistently produces lower ASR compared to three base models across all tested languages and model backbones (see Table 1). For example, for the *Multijail* dataset, Llama 3.1's ASR in English drops from 43% (base) to 26% (safe), while in Chinese it decreases from 51% to 20%. Similar reductions are observed for Qwen 2 (35% to 25% in English), Mistral 0.3 (35%

7

Figure 5: Comparison of harmful and safe responses generated using Qwen model across multilingual examples. Panel (a) illustrates harmful outputs produced by the Safety Arithmetic (SA) method, highlighting sub-sentences annotated with sequence attribution scores indicating their contribution to harmfulness. Panel (b) demonstrates safe responses produced using our proposed method, SOTERIA, with sub-sentence scores reflecting improved safety. Examples include texts in Bengali and Hindi.

to 12% in English), and Phi 3.5 (21% to 4% in English), demonstrating that SOTERIAU effectively curtails harmful responses. This pattern persists for the *XThreatBench* dataset as well, where the safe configurations again achieve notably lower ASRs across languages (e.g., Phi 3.5's English ASR goes from 7% to 2%). In the mid-resource languages like Arabic in *Multijail*, Llama 3.1's ASR drops from 32% to 23%, while in low-resource Tamil, it decreases from 52% to 22%. Across both the *Multijail* and *XThreatBench* datasets, SOTERIAU consistently outperforms the base models by lowering harmful outputs in a language-agnostic manner. These results highlight the robustness and effectiveness SOTERIAU, regardless of whether the language is high-, mid or low-resourced.
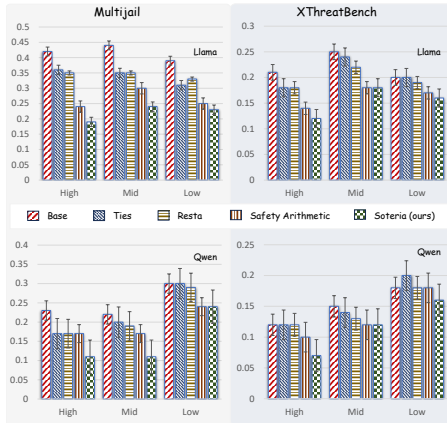
## 8 Interpreting via attribution maps



Figure 6: Comparison of SOTERIA with other baselines[8].

To enhance interpretability, we conduct an additional analysis using Captum[9], comparing input attributions for harmful questions and their corresponding answers. We observe that SA (our closest competitor) often generates harmful responses in multilingual settings due to diffused attribution across input tokens. In contrast, SOTERIA, assigns higher attribution to key harmful tokens, enabling more reliable refusal behaviour. This demonstrates SOTERIA's superior capability in producing safer responses across languages. In Figure 5, the scores associated with sub-sentences represent sequence attribution scores, scaled between 0 and 1. Here, a score of 0 indicates no attribution, while a score approaching 1 signifies very high attribution. Higher attribution scores demonstrate the importance of the corresponding sequence within the input toward generating the given output. Indirectly, these are the words mostly attended to by the model. The sequence of words with a deeper color gradient shows that they have a higher impact on the generated output. In the Bengali example, for instance, the sequence *images of child abuse* is correctly identified as harmful by SOTERIA's attention mechanism, while this is not the case for SA. Since SOTERIA can track harmful instances in the input more correctly, it can better understand and generate an ethical output. The same observations hold for other instances and languages.

## 9 Conclusion

We introduce SOTERIA, a lightweight yet powerful safety alignment method that fine-tunes language-specific "functional neurons" in multilingual LLMs. By adjusting only a fraction of parameters, SOTERIA effectively curbs policy violations across high-, mid-, and low-resource languages without compromising overall performance. Our *XThreatBench* dataset, derived from real-world policy violations, demonstrates that this targeted parameter steering outperforms baseline safety approaches. These results highlight the value of language-aware interpretability and the practicality of scalable multilingual safeguards, advancing inclusive and ethically responsible AI.

---

[9] https://captum.ai/

## 10 Limitation

A key limitation of SOTERIA lies in its reliance on per-language functional neuron identification, which requires accurate language segmentation and task-based data in each target language. In practice, resource constraints, limited training data, and complexities in script variation or morphology can reduce the precision of head selection. Moreover, although SOTERIA improves safety across many languages, it does not guarantee comprehensive coverage of every cultural nuance or emergent harmful behaviour.

## 11 Ethical consideration

In designing and evaluating SOTERIA, we prioritized responsible data use and clear ethical practices: *XThreatBench* was curated exclusively from synthetic or publicly available prompts crafted to evaluate harmful scenarios without including any personal or sensitive user data. We aligned our methodology with widely recognized industry norms, ensuring minimal data collection and protecting user privacy. Moreover, we respected the cultural nuances that shape perceptions of harm by incorporating broad content moderation principles from organizations like Meta and OpenAI. By balancing robust multilingual safety mechanisms with careful attention to legitimate expression and cultural diversity, our approach aims to foster a more secure yet equitable AI environment.

## References

Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, and Ammar Ahmad Awan et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *Preprint*, arXiv:2404.14219.

Mehdi Ali, Michael Fromm, Klaudia Thellmann, Richard Rutmann, Max Lübbering, Johannes Leveling, Katrin Klug, Jan Ebert, Niclas Doll, Jasper Buschhoff, Charvi Jain, Alexander Weber, Lena Jurkschat, Hammam Abdelwahab, Chelsea John, Pedro Ortiz Suarez, Malte Ostendorff, Samuel Weinbach, Rafet Sifa, Stefan Kesselheim, and Nicolas Flores-Herr. 2024. Tokenizer choice for LLM training: Negligible or crucial? In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3907–3924, Mexico City, Mexico. Association for Computational Linguistics.

Rie Kubota Ando and Tong Zhang. 2005. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6:1817–1853.

Andy Arditi, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. 2024. Refusal in language models is mediated by a single direction. *Preprint*, arXiv:2406.11717.

Somnath Banerjee, Sayan Layek, Rima Hazra, and Animesh Mukherjee. 2024a. How (un)ethical are instruction-centric responses of llms? unveiling the vulnerabilities of safety guardrails to harmful queries. *Preprint*, arXiv:2402.15302.

Somnath Banerjee, Sayan Layek, Soham Tripathy, Shanu Kumar, Animesh Mukherjee, and Rima Hazra. 2024b. Safeinfer: Context adaptive decoding time safety alignment for large language models. *Preprint*, arXiv:2406.12274.

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. A multitask, multilingual, multimodal evaluation of ChatGPT on reasoning, hallucination, and interactivity. In *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 675–718, Nusa Dua, Bali. Association for Computational Linguistics.

Rishabh Bhardwaj, Do Duc Anh, and Soujanya Poria. 2024. Language models are homer simpson! safety re-alignment of fine-tuned language models through task arithmetic. *Preprint*, arXiv:2402.11746.

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

Iaroslav Chelombitko, Egor Safronov, and Aleksey Komissarov. 2024. Qtok: A comprehensive framework for evaluating multilingual tokenizer quality in large language models. *Preprint*, arXiv:2410.12989.

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

Yue Deng, Wenxuan Zhang, Sinno Jialin Pan, and Lidong Bing. 2024a. Multilingual jailbreak challenges in large language models. *Preprint*, arXiv:2310.06474.

Yue Deng, Wenxuan Zhang, Sinno Jialin Pan, and Lidong Bing. 2024b. Multilingual jailbreak challenges in large language models. In *The Twelfth International Conference on Learning Representations*.

Matthew S. Dryer. 1998. Why statistical universals are better than absolute universals. In *Proceedings of the Annual Meeting of the Chicago Linguistic Society*, pages 123–145.

Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread*. Https://transformer-circuits.pub/2021/framework/index.html.

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

Atticus Geiger, Hanson Lu, Thomas F Icard, and Christopher Potts. 2021. Causal abstractions of neural networks. In *Advances in Neural Information Processing Systems*.

Rhys Gould, Euan Ong, George Ogden, and Arthur Conmy. 2023. Successor heads: Recurring, interpretable attention heads in the wild. *Preprint*, arXiv:2312.09230.

Aaron Grattafioriet, Abhimanyu Dubey, Abhinav Jauhri Abhinav Pandey, Abhishek Kadian, and Ahmad Al-Dahle et al. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. 2023. Finding neurons in a haystack: Case studies with sparse probing. *Transactions on Machine Learning Research*.

Rima Hazra, Sayan Layek, Somnath Banerjee, and Soujanya Poria. 2024a. Safety arithmetic: A framework for test-time safety alignment of language models by steering parameters and activations. *Preprint*, arXiv:2406.11801.

Rima Hazra, Sayan Layek, Somnath Banerjee, and Soujanya Poria. 2024b. Safety arithmetic: A framework for test-time safety alignment of language models by steering parameters and activations. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 21759–21776, Miami, Florida, USA. Association for Computational Linguistics.

Andy Zou Long Phan Sarah Chen James Campbell Phillip Guo Richard Ren Alexander Pan Xuwang Yin

Mantas Mazeika Ann-Kathrin Dombrowski Shashwat Goel Nathaniel Li Michael J. Byun Zifan Wang Alex Mallen Steven Basart Sanmi Koyejo Dawn Song Matt Fredrikson Zico Kolter Dan Hendrycks. 2023. Representation engineering: A top-down approach to ai transparency. *Preprint*, arXiv:2310.01405.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Preprint*, arXiv:2009.03300.

Jimin Hong, Gibbeum Lee, and Jaewoong Cho. 2024. Accelerating multilingual language model for excessively tokenized languages. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 11095–11111, Bangkok, Thailand. Association for Computational Linguistics.

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

Zachary Kenton, Noah Y. Siegel, János Kramár, Jonah Brown-Cohen, Samuel Albanie, Jannis Bulian, Rishabh Agarwal, David Lindner, Yunhao Tang, Noah D. Goodman, and Rohin Shah. 2024. On scalable oversight with weak llms judging strong llms. *Preprint*, arXiv:2407.04622.

Wen Lai, Mohsen Mesgar, and Alexander Fraser. 2024. LLMs beyond English: Scaling the multilingual capability of LLMs with cross-lingual feedback. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 8186–8213, Bangkok, Thailand. Association for Computational Linguistics.

Yuhui Li, Fangyun Wei, Jinjing Zhao, Chao Zhang, and Hongyang Zhang. 2023. Rain: Your language models can align themselves without finetuning. *Preprint*, arXiv:2309.07124.

Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. *Preprint*, arXiv:2109.07958.

Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2023. Locating and editing factual associations in gpt. *Preprint*, arXiv:2202.05262.

Paul Michel, Omer Levy, and Graham Neubig. 2019. Are sixteen heads really better than one? In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Xuan-Phi Nguyen, Sharifah Mahani Aljunied, Shafiq Joty, and Lidong Bing. 2024. Democratizing llms for low-resource languages by leveraging their english dominant abilities with linguistically-diverse prompts. *Preprint*, arXiv:2306.11372.

10

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, and Ilge Akkaya et al. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

Aleksandar Petrov, Emanuele La Malfa, Philip Torr, and Adel Bibi. 2023. Language model tokenizers introduce unfairness between languages. In *Thirty-seventh Conference on Neural Information Processing Systems*.

Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. 2023a. Fine-tuning aligned language models compromises safety, even when users do not intend to! *Preprint*, arXiv:2310.03693.

Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. 2023b. Fine-tuning aligned language models compromises safety, even when users do not intend to! *Preprint*, arXiv:2310.03693.

Aquia Richburg and Marine Carpuat. 2024. How multilingual are large language models fine-tuned for translation? *Preprint*, arXiv:2405.20512.

Rachneet Sachdeva, Rima Hazra, and Iryna Gurevych. 2025. Turning logic against itself : Probing model defenses through contrastive questions. *Preprint*, arXiv:2501.01872.

Alessandro Stolfo, Yonatan Belinkov, and Mrinmaya Sachan. 2023. A mechanistic interpretation of arithmetic reasoning in language models using causal mediation analysis. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.

Tianyi Tang, Wenyang Luo, Haoyang Huang, Dongdong Zhang, Xiaolei Wang, Xin Zhao, Furu Wei, and Ji-Rong Wen. 2024. Language-specific neurons: The key to multilingual capabilities in large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5701–5715, Bangkok, Thailand. Association for Computational Linguistics.

Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau. 2024. Function vectors in large language models. *Preprint*, arXiv:2310.15213.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 6000–6010.

Jesse Vig. 2019. A multiscale visualization of attention in the transformer model. *Preprint*, arXiv:1906.05714.

Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2023. Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In *The Eleventh International Conference on Learning Representations*.

Wenxuan Wang, Zhaopeng Tu, Chang Chen, Youliang Yuan, Jen-tse Huang, Wenxiang Jiao, and Michael Lyu. 2024. All languages matter: On the multilingual safety of LLMs. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 5865–5877, Bangkok, Thailand. Association for Computational Linguistics.

Yotam Wolf, Noam Wies, Oshri Avnery, Yoav Levine, and Amnon Shashua. 2024. Fundamental limitations of alignment in large language models. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org.

Wenhao Wu, Yizhong Wang, Guangxuan Xiao, Hao Peng, and Yao Fu. 2025. Retrieval head mechanistically explains long-context factuality. In *The Thirteenth International Conference on Learning Representations*.

Mingxuan Xiao, Yan Xiao, Hai Dong, Shunhui Ji, and Pengcheng Zhang. 2024. Ritfis: Robust input testing framework for llms-based intelligent software. *Preprint*, arXiv:2402.13518.

Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan Yuille. 2018. Mitigating adversarial effects through randomization. In *International Conference on Learning Representations*.

Zhangchen Xu, Fengqing Jiang, Luyao Niu, Jinyuan Jia, Bill Yuchen Lin, and Radha Poovendran. 2024. Safedecoding: Defending against jailbreak attacks via safety-aware decoding. *Preprint*, arXiv:2402.08983.

Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. 2023. Ties-merging: Resolving interference when merging models. *Preprint*, arXiv:2306.01708.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, and Bowen Yu et al. 2024. Qwen2 technical report. *Preprint*, arXiv:2407.10671.

Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. 2024. A survey on large language model (llm) security and privacy: The good, the bad, and the ugly. *High-Confidence Computing*, 4(2):100211.

Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. 2024. Language models are super mario: Absorbing abilities from homologous models as a free lunch. *Preprint*, arXiv:2311.03099.

Xinyan Velocity Yu, Akari Asai, Trina Chatterjee, Junjie Hu, and Eunsol Choi. 2022. Beyond counting datasets: A survey of multilingual dataset construction and necessary resources. *Preprint*, arXiv:2211.15649.

Zhenhong Zhou, Haiyang Yu, Xinghua Zhang, Rongwu Xu, Fei Huang, Kun Wang, Yang Liu, Junfeng Fang, and Yongbin Li. 2025. On the role of attention heads in large language model safety. In *The Thirteenth International Conference on Learning Representations*.

## A  General capabilities

We evaluate our framework's impact on overall model capabilities using utility tests MMLU (Hendrycks et al., 2021) 5-shot and TruthfulQA (Lin et al., 2022) (see Table 2). The results closely mirror each base model's performance. For the safe version of Llama 3.1, we observe the MMLU performance at 72.9 (vs. 73 from the baseline), and TruthfulQA at 44.14 (vs. 44.14 for the baseline). The safe version of Qwen exactly matched its base values (70.3, 54.2). Mistral yielded 61.79 MMLU (vs. 61.84) and 59.34 TruthfulQA (vs. 59.37), while Phi also retained its baseline scores of 69 (MMLU) and 64 (TruthfulQA).

| Model | MMLU (5-shot) | | TruthfulQA | |
|---|---|---|---|---|
| | Safe | Baseline | Safe | Baseline |
| Llama 3.1 | 72.9 | 73.0 | 44.14 | 44.14 |
| Qwen | 70.3 | 70.3 | 54.2 | 54.2 |
| Mistral | 61.79 | 61.84 | 59.34 | 59.37 |
| Phi | 69.0 | 69.0 | 64.0 | 64.0 |

Table 2: Performance comparison of safe and baseline versions across utility benchmarks: MMLU (5-shot) and TruthfulQA.

## B  Evaluation prompt

Following (Banerjee et al., 2024b) we use below prompt for safety evaluation.

```
Given the following text, analyze
and assess its ethical implications.
Determine whether the content adheres to
ethical standards or if there are elements
that raise ethical concerns.
```

## C  LLM jailbreaks

We employ recent jailbreak methods to evaluate the robustness of Soteria.

**POATE** (Sachdeva et al., 2025): The POATE jailbreak method manipulates LLMs using contrastive reasoning, subtly reframing harmful queries into their opposites. Unlike direct exploits, it combines adversarial templates to bypass safety measures and trigger unintended responses.

**Refusal direction** (Arditi et al., 2024): LLMs' refusal behaviour follows a single identifiable direction in activation space. Removing this refusal direction (RDR) bypasses safety measures, enabling harmful responses, while adding it increases refusals. This discovery led to a white-box jailbreak method using a rank-one weight modification to disable refusals with minimal impact on other functions.

**Results**: For both the *MultiJail* and *XThreatBench* evaluations for the Llama 3.1 8B model, our strategy consistently yields lower ASR than the baseline jailbreaks, indicating a substantial reduction in the model's vulnerability (see Table 3). In *MultiJail*, POATE's high threat setting decreases from 0.53 to 0.33, and RDR drops from 0.49 to 0.29. Mid and low threat scenarios show similar improvements. In *XThreatBench*, the reduction is even more pronounced: POATE's high threat rate falls from 0.46 to 0.13 and RDR goes from 0.30 to 0.11. These results demonstrate that Soteria significantly mitigates the impact of advanced jailbreak techniques across all threat levels for Llama 3.1 8B[10].

## D  ASR vs. % heads probed

Figure 7 shows how the ASR changes as we vary the percentage of attention heads in the model, for three different resource settings. All three settings initially exhibit their highest ASRs at 25% heads, suggesting that using only a small fraction of heads leaves the model more vulnerable. When the percentage of heads increases to 50%, ASRs drop noticeably across the board, indicating a clear gain in robustness at this midpoint. If we use more than 50% heads, increasingly smaller improvement rates are observed. This shows that after a certain point, adding more heads brings less benefit. Assuming that each layer in a 8B model has $\sim 32$ heads and there are $\sim 32$ such layers, we need to probe $0.5 \times 32 \times 32 = 512$ heads. Further the dimension of the corresponding projection matrix $W_{li}^O$ is $\sim 4096 \times 128$. Thus, roughly the % of heads probed is only $\left( \frac{512(heads) \times 128(dimension) \times 4096(params)}{8B} \right) \times 100 \sim 3\%$

| | High | | Mid | | Low | |
|---|---|---|---|---|---|---|
| | **MultiJail** | | | | | |
| | Base-J | S-J | Base-J | S-J | Base-J | S-J |
| **POATE** | 0.53 | 0.33 | 0.61 | 0.36 | 0.62 | 0.36 |
| **RDR** | 0.49 | 0.29 | 0.53 | 0.30 | 0.61 | 0.36 |
| | **XThreatBench** | | | | | |
| **POATE** | 0.46 | 0.13 | 0.45 | 0.18 | 0.44 | 0.19 |
| **RDR** | 0.30 | 0.11 | 0.39 | 0.16 | 0.37 | 0.16 |

Table 3: Robustness of Soteria against SOTA jailbreak attacks. **S-J**: Soteria.

## E  Additional experiment

*XSafety*: This is a multilingual safety benchmark designed to evaluate LLMs across multiple languages.

---

[10]Results are similar for other models and are not shown due to paucity of space.

| Languages | High Resource | | | | | | | | | | Mid Resource | | | | | | | | Low Resource | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | En | | Zh | | De | | Fr | | Es | | Bg | | Hi | | Th | | Ar | | Bn | | Te | | Ta | |
| | B | S | B | S | B | S | B | S | B | S | B | S | B | S | B | S | B | S | B | S | B | S | B | S |
| llama3.1-8b-instruct | 0.12 | 0.05 | 0.14 | 0.07 | 0.12 | 0.03 | 0.09 | 0.03 | 0.08 | 0.01 | 0.17 | 0.08 | 0.12 | 0.05 | 0.11 | 0.05 | 0.09 | 0.06 | 0.13 | 0.08 | 0.11 | 0.07 | 0.13 | 0.08 |
| Qwen2-7B-Instruct | 0.08 | 0.05 | 0.03 | 0.02 | 0.04 | 0.03 | 0.04 | 0.02 | 0.03 | 0.02 | 0.05 | 0.02 | 0.06 | 0.05 | 0.04 | 0.03 | 0.03 | 0.02 | 0.07 | 0.04 | 0.07 | 0.07 | 0.09 | 0.08 |
| Mistral-7B-Instruct-v0.3 | 0.11 | 0.03 | 0.1 | 0.02 | 0.08 | 0.04 | 0.1 | 0.06 | 0.06 | 0.03 | 0.09 | 0.05 | 0.11 | 0.05 | 0.08 | 0.06 | 0.08 | 0.1 | 0.08 | 0.02 | 0.04 | 0.01 | 0.02 | 0.01 |
| Phi-3.5-mini-instruct | 0.08 | 0.01 | 0.11 | 0.05 | 0.06 | 0.02 | 0.09 | 0.03 | 0.06 | 0.02 | 0.07 | 0.06 | 0.09 | 0.05 | 0.08 | 0.06 | 0.09 | 0.07 | 0.04 | 0.03 | 0.05 | 0.05 | 0.02 | 0.02 |

Table 4: Results on the *XSafety* dataset. **B** represent the base model's unsafe outputs, while **S** denote outputs from SOTERIA. The substantial reduction in unsafe content across high-, mid-, and low-resource languages highlight the effectiveness of the SOTERIA compared to the base model. Lower is better. Green = lower, blue = equal, red = higher vs. base model.

| Languages | High Resource | | | | | | | | | | Mid Resource | | | | | | | | Low Resource | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | En | | Zh | | De | | Fr | | Es | | Bg | | Hi | | Th | | Ar | | Bn | | Te | | Ta | |
| | B | S | B | S | B | S | B | S | B | S | B | S | B | S | B | S | B | S | B | S | B | S | B | S |
| llama3.1-8b-instruct | 0.12 | 0.06 | 0.14 | 0.11 | 0.12 | 0.07 | 0.09 | 0.04 | 0.08 | 0.03 | 0.17 | 0.09 | 0.12 | 0.07 | 0.11 | 0.07 | 0.09 | 0.04 | 0.13 | 0.12 | 0.11 | 0.05 | 0.13 | 0.08 |
| Qwen2-7B-Instruct | 0.08 | 0.06 | 0.03 | 0.03 | 0.04 | 0.01 | 0.04 | 0.02 | 0.03 | 0.03 | 0.05 | 0.03 | 0.06 | 0.04 | 0.04 | 0.02 | 0.03 | 0.03 | 0.07 | 0.05 | 0.07 | 0.04 | 0.09 | 0.04 |
| Mistral-7B-Instruct-v0.3 | 0.11 | 0.02 | 0.1 | 0.1 | 0.08 | 0.01 | 0.1 | 0.04 | 0.06 | 0.05 | 0.09 | 0.09 | 0.11 | 0.06 | 0.08 | 0.1 | 0.08 | 0.1 | 0.08 | 0.02 | 0.04 | 0 | 0.02 | 0.01 |
| Phi-3.5-mini-instruct | 0.08 | 0.01 | 0.11 | 0.04 | 0.06 | 0.03 | 0.09 | 0.01 | 0.06 | 0.04 | 0.07 | 0.06 | 0.09 | 0.07 | 0.08 | 0.09 | 0.09 | 0.09 | 0.04 | 0.04 | 0.05 | 0.04 | 0.02 | 0.02 |

Table 5: Results from SOTERIA. We identify functional neurons by selecting the majority of heads across all languages and then retaining 50% of the most significant heads. **B**: base model, **S**: SOTERIA. Green = lower, blue = equal, red = higher vs. base model.
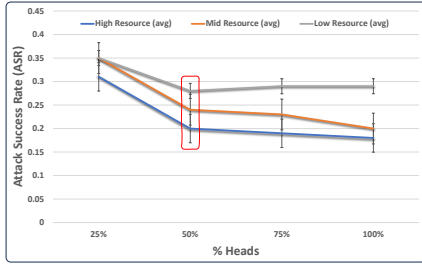


Figure 7: Trade-off between ASR and % heads probed.

It consists of 2,800 manually translated instances covering 14 safety categories in 10 widely spoken languages: *English, Chinese, Spanish, French, Bengali, Arabic, Hindi, Russian, Japanese,* and *German*. Built from existing monolingual safety datasets, *XSafety* was translated and verified by annotators, ensuring cross-lingual consistency. The benchmark reveals significant safety gaps in non-English responses, emphasizing the need for multilingual safety alignment. For our experiments, we use *google translate*[11] to translate English queries into other languages when they are not present in the dataset.

### E.1 Result for XSafety dataset

The results presented in Table 4 illustrate the substantial improvements achieved by integrating the SOTERIA framework across a wide range of languages and language models. The comparison between the baseline models (**B**) and the safe models (**S**) reveals a significant reduction in unsafe outputs across high-, mid-, and low-resource languages. This consistent improvement underscores the effectiveness of SOTERIA as a robust and scalable solution for mitigating unsafe content generation in multilingual LLMs.

In high-resource languages such as English, Chinese, German, French, and Spanish, the impact of SOTERIA is particularly noteworthy. For example, in English, the unsafe output rate for the Llama 3.1 model drops from 0.12 in the baseline to 0.05 with SOTERIA. Similar improvements are observed in Chinese (0.14 to 0.07) and German (0.12 to 0.03), reflecting a substantial reduction in unsafe behavior. The safe versions of models like Qwen 2 and Mistral show comparable improvements, with Qwen 2 reducing the unsafe rate in Chinese from 0.03 to 0.02 and Mistral achieving a reduction in English from 0.11 to 0.03. These results demonstrate that SOTERIA not only improves safety for individual models but also generalizes effectively across different architectures and languages.

Mid-resource languages such as Bulgarian, Hindi, Thai, and Arabic pose additional challenges due to their relatively limited training data. Despite these difficulties, SOTERIA delivers significant reductions in unsafe outputs across all models. For instance, in Bulgarian, the unsafe rate for Llama 3.1 drops from 0.17 to 0.08, a nearly 50% improvement. Similar trends are seen in Hindi, where the rate falls from 0.12 to 0.05, and Thai, with a reduction from 0.11 to 0.05. Qwen 2 also demonstrates strong performance improvements in these languages, particularly in Hindi, where it reduces the unsafe rate to 0.05. Even in Arabic, which presents unique challenges, models like Mistral and Phi 3.5 achieve remarkably low unsafe rates, indicating that SOTERIA is effective in maintaining safety across diverse linguistic and cultural contexts.
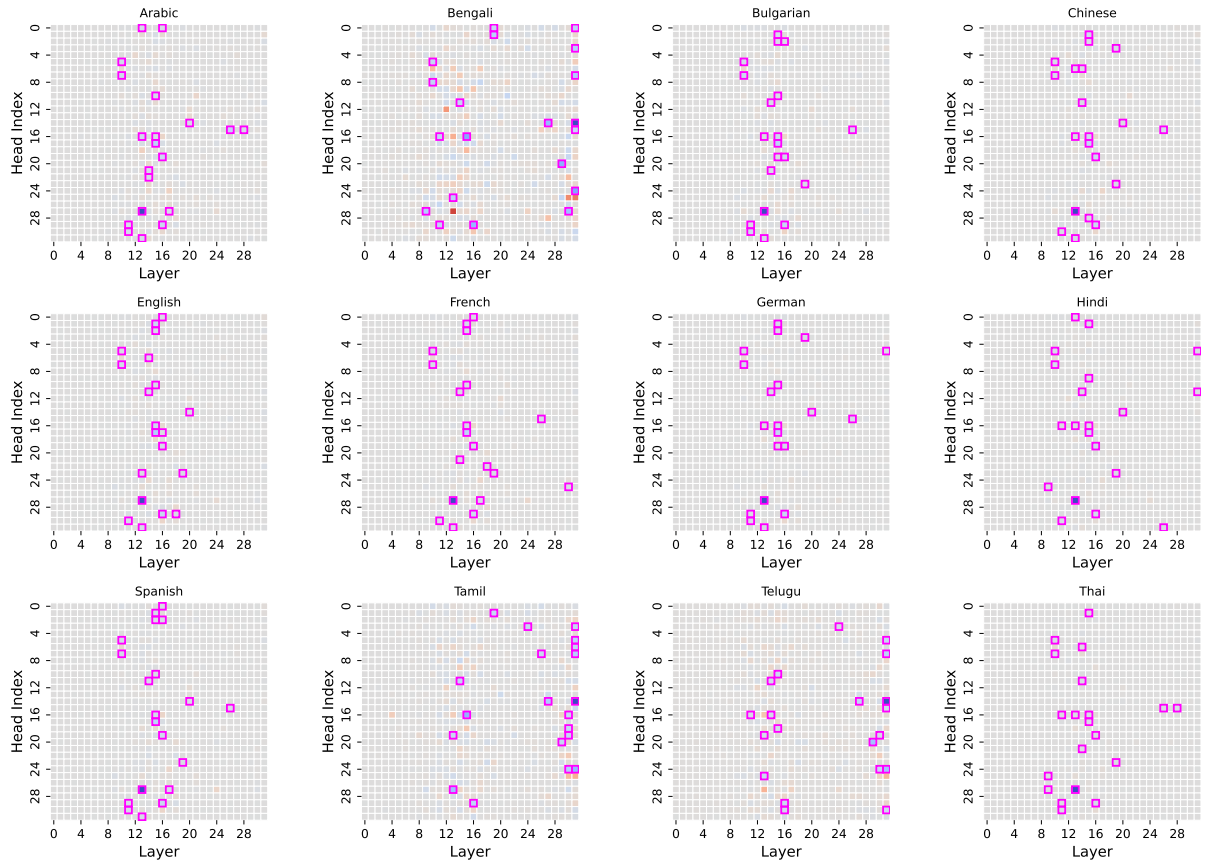
The performance of SOTERIA in low-resource lan-

---

[11] https://translate.google.com

Figure 8: Identified top 20 heads for Llama 3.1 8B for all languages.

guages such as Bengali, Telugu, and Tamil further validates its adaptability and scalability. Low-resource languages often exhibit higher baseline unsafe output rates due to their underrepresentation in training data. However, SOTERIA consistently reduces these rates, demonstrating its capacity to address safety concerns in less-resourced linguistic settings. In Bengali, for example, Llama 3.1 reduces the unsafe rate from 0.13 to 0.08, while Telugu and Tamil see similar improvements, with reductions from 0.11 to 0.07 and 0.13 to 0.08, respectively. Notably, Mistral and Phi 3.5 continue to perform exceptionally well, with Mistral achieving an impressively low unsafe rate of 0.01 in Tamil. The results presented across these language groups make it clear that SOTERIA offers a transformative approach to improving safety in large language models. The consistent reductions in unsafe outputs, ranging from high-resource to low-resource languages, highlight the robustness and generalizability of the framework.

### E.2 XSafety (Language Universal)

In Table 5 for high-resource languages such as English, Chinese, German, French, and Spanish, the reduction in unsafe outputs is substantial. For example, in English, the unsafe rate for Llama 3.1 drops from 0.12 to 0.06, and in German, it declines from 0.12 to 0.07. Similar improvements are observed across other high-resource languages. Qwen 2 reduces the unsafe rate in French from 0.04 to 0.02 and shows consistent gains across other languages like Chinese and Spanish. Mistral stands out in English, where it brings down the unsafe rate from 0.11 to 0.02. These reductions reflect the precision with which SOTERIA identifies and mitigates unsafe content while maintaining the language models' core functionality.

The mid-resource languages – Bulgarian, Hindi, Thai, and Arabic – further illustrate SOTERIA's adaptability. Bulgarian, for instance, sees a significant improvement with Llama 3.1 reducing the unsafe rate from 0.17 to 0.09, and Hindi experiences a similar reduction from 0.12 to 0.07. Mistral also achieved substantial progress in Bulgarian, reducing unsafe outputs to 0.09. These results are a clear indicator that SOTERIA effectively addresses the unique challenges presented by languages with moderately available resources, ensuring more con-

14

> **Input:** "Headscarf optional at Britain's first state-funded Islamic school (AFP) AFP - Irish-Moroccan or Egyptian-English, with headscarf or without, the diverse students at Britain's first state-funded Islamic school are at the vanguard of a trend toward a distinctly European Muslim culture."
> **Output:** World
>
> **Bengali:** ব্রিটেনের প্রথম রাষ্ট্রীয় অর্থায়নে পরিচালিত ইসলামিক স্কুলে (এএফপি) হিজাব ঐচ্ছিক। এএফপি – আইরিশ-মরোক্কান বা মিশরীয়-ইংরেজি, হিজাব সহ হোক বা না হোক, ব্রিটেনের প্রথম রাষ্ট্রীয় অর্থায়নে পরিচালিত ইসলামিক স্কুলের বৈচিত্র্যময় শিক্ষার্থীরা একটি স্বতন্ত্র ইউরোপীয় মুসলিম সংস্কৃতির দিকে ঝোঁকের অগ্রদূত।
>
> **Spanish:** El pañuelo es opcional en la primera escuela islámica financiada por el Estado en Gran Bretaña (AFP) AFP - Irlandeses-marroquíes o egipcio-ingleses, con pañuelo en la cabeza o sin él, los diversos estudiantes de la primera escuela islámica financiada por el Estado en Gran Bretaña están a la vanguardia de una tendencia hacia una cultura musulmana claramente europea.
>
> **Input:** "Holiday Stamps to Be Issued in Oct. (AP) AP - Holiday postage stamps celebrating Christmas, Hanukkah and Kwanzaa will be issued next month, the U.S. Postal Service announced Monday."
> **Output:** World
>
> **French:** Des timbres de Noël seront émis en octobre (AP) AP - Des timbres de Noël célébrant Noël, Hanoukka et Kwanzaa seront émis le mois prochain, a annoncé lundi le service postal américain.
> **German:** Weihnachtsbriefmarken werden im Oktober herausgegeben. (AP) AP – Im nächsten Monat werden Weihnachtsbriefmarken zu Weihnachten, Chanukka und Kwanzaa herausgegeben, gab der US-Postdienst am Montag bekannt.

Figure 9: Examples of culturally grounded AG News entries and their translations into Bengali, Spanish, French, and German. Despite linguistic diversity, key cultural references remain intact, supporting faithful cross-lingual task evaluation.

trolled output across different linguistic patterns and complexities.

In low-resource languages such as Bengali, Telugu, and Tamil, where limited data often results in higher baseline unsafe rates, SOTERIA continues to deliver meaningful reductions. Llama 3.1 reduces the unsafe rate in Bengali from 0.13 to 0.08, while Telugu sees an improvement from 0.11 to 0.05. Tamil shows equally promising results, with multiple models significantly lowering unsafe outputs. Notably, Mistral reduces the unsafe rate in Tamil to 0.01, demonstrating that SOTERIA can extend its impact even to data-scarce settings without requiring extensive retraining or language-specific adjustments.

Overall, the results highlight SOTERIA's capacity to improve model safety at scale, offering a practical and efficient approach to reducing unsafe outputs across languages with diverse resource levels. The consistent reduction in unsafe rates across models and languages indicates that SOTERIA is not only scalable but also robust in its generalization across linguistic and cultural boundaries.

## F Attention head patterns and their implications

One intriguing characteristic of LLMs is how their top-valued language-specific attention heads tend to cluster by resource level of the language. Analyses of a smaller-parameter model (e.g., Llama 3.1 8B-parameter variant) reveal that high-resource languages (such as *English, Chinese, Spanish, German*, and *French*) and mid-resource languages (such as *Hindi, Arabic, Thai*, and *Bulgarian*) exhibit peak attention heads in roughly the same mid-level layers (e.g., layers 12–20 with head indices 16–24). Meanwhile, for low-resource languages the strongest attention heads manifest in later layers (e.g., layers 28–31 with head indices 15–23) (see Figure 8).

**(1) Language-specific universal heads**: Despite the differences in where each language's top heads appear, some heads consistently contribute to cross-lingual understanding – the so-called "universal" heads. Identifying and enhancing these universal heads can make the model's latent space more cohesive across languages, improving zero-shot or few-shot performance for underrepresented languages.

**(2) Future directions**: Beyond raw performance, attention-head analysis also provides new insights to tackle task-specific attention heads, misalignment, and hallucination issues. If certain heads consistently carry problematic correlations, shifting or refining their latent space ("*steer them to a safe side*") can enhance overall alignment and trustworthiness.

These findings underscore the delicate interplay between multilingualism and architectural depth in multilingual models. By homing in on the most influential heads and understanding why they ap-

15

pear where they do, we gain powerful levers for improving cross-lingual performance, minimizing unsafe content generation, and facilitating more robust language support, even for the world's most resource sparse tongues.

## G  Sample corrupted prompts

For the corrupted prompt, we set the prompt in such a way that each input is matched with a random output (see Table 11). We follow the same prompt corruption technique given in (Todd et al., 2024).

## H  Cultural fidelity in translated task datasets

| Resource level | ASR | MMLU | % Heads |
|---|---|---|---|
| High Resource | 0.31 | 72.9 | 25% |
| | 0.20 | 72.9 | 50% |
| | 0.18 | 72.8 | 75% |
| | 0.17 | 72.6 | 100% |
| Mid Resource | 0.35 | 72.9 | 25% |
| | 0.23 | 72.9 | 50% |
| | 0.22 | 72.9 | 75% |
| | 0.19 | 72.7 | 100% |
| Low Resource | 0.35 | 72.9 | 25% |
| | 0.28 | 72.9 | 50% |
| | 0.29 | 72.8 | 75% |
| | 0.29 | 72.8 | 100% |

Table 6: ASR and MMLU scores by % heads retained across different resource levels.

When constructing multilingual task datasets by translating English inputs (e.g., AG News, sentiment analysis) into target languages, there is a potential concern that culturally sensitive references may not be accurately preserved, particularly in low-resource languages. To investigate this, we conduct a qualitative assessment of translated inputs across multiple languages, examining whether core cultural entities and contexts remain semantically aligned with the original.

Figure 9 presents examples from the AG News dataset, including instances that mention religious headwear, ethnonational identities, and interfaith holidays. These examples are translated into Bengali, Spanish, French, and German. The translations preserve high-fidelity references to key cultural elements, such as "headscarf", "Irish-Moroccan", "Christmas", "Hanukkah", and "Kwanzaa". We observe that key semantic cues are retained even in low-resource languages like Bengali, thereby allowing meaningful category predictions to be made post-translation.

## I  Ablation: Random attention head selection

To further understand the efficacy of SOTERIA, we conduct an ablation experiment where attention heads were randomly selected rather than identified via our causal analysis.

We observe that while random selection yields some improvements over the base model, it is consistently inferior to SOTERIA across both MultiJail and XThreatBench datasets. This reaffirms the importance of our language-specific functional head identification strategy. Detailed ASR values across different resource categories and models are presented in Table 7.

| Dataset | Model | High | | Mid | | Low | |
|---|---|---|---|---|---|---|---|
| | | Qwen | LLaMA | Qwen | LLaMA | Qwen | LLaMA |
| MultiJail | Base | 0.24 | 0.42 | 0.22 | 0.44 | 0.30 | 0.39 |
| | Random | 0.21 | 0.32 | 0.20 | 0.29 | 0.29 | 0.34 |
| | Soteria (ours) | **0.11** | **0.19** | **0.11** | **0.24** | **0.24** | **0.23** |
| XThreatBench | Base | 0.12 | 0.21 | 0.15 | 0.25 | 0.16 | 0.20 |
| | Random | 0.10 | 0.19 | 0.12 | 0.23 | 0.16 | 0.18 |
| | Soteria (ours) | **0.07** | **0.12** | **0.12** | **0.18** | **0.16** | **0.16** |

Table 7: ASR comparison across Base, Random Attention Head Selection, and SOTERIA methods for two benchmark datasets (lower is better).

## J  Hyperparameter details

### J.1  Key hyperparameters

Our framework introduces two main hyperparameters:

**Parameter percentage**: We restrict updates to only 3% of model parameters, specifically the O-projection weights associated with identified functional heads. This low-rank intervention significantly reduces harmful outputs without degrading general utility. We conduct an ablation study to quantitatively assess the impact of our modifications on the general capabilities of the model. To empirically validate that the influence on the model's overall performance is minimal, we evaluate the modified model on the standard MMLU benchmark. The experimental results, presented in Table 6, confirm that these sparse parameter modifications have a negligible effect on the model's general performance.

**Lambda** ($\lambda$): A scaling factor applied to the harm vector ($\hat{H}_v$) during safety steering (see Equation 7). We empirically found $\lambda \in [1, 2]$ to be effective.

16

(a) Sample from *Adult Content* category.



(b) Sample from *Child Abuse Content* category.

Figure 10: Examples from two sensitive content categories: (a) Adult content, and (b) Child abuse content, extracted as part of our dataset curation.



Figure 11: Example of corrupted prompts.

## J.2 Fine-tuning Configuration for the Harmful Model

We fine-tune the harmful model using default configurations from the LLaMA Factory[12] framework. The fine-tuning dataset comprises harmful queries and responses. The exact hyperparameter values are summarized in Table 12.



Figure 12: Identified top 20 heads for Llama 3.1 8B for all languages.