

# Unlocking LLM Safeguards for Low-Resource Languages via Reasoning and Alignment with Minimal Training Data

Zhuowei Chen, Bowei Zhang, Nankai Lin, Tian Hou, Lianxi Wang

Guangdong University of Foreign Studies & University of Pittsburgh

## CHALLENGES

### LLM Safeguards are Essential:

Recent LLM advances increase the risk of malicious requests.  
(e.g., jailbreaking, harmful content generation).

### Current Limitations:

1. Lack Interpretability (Black-box problem).
2. Poor Performance on Low-Resource Languages (LRLs).
3. Require Significantly Large Training Datasets.

## METHODOLOGY

### Reasoning-based Multilingual Safeguards

- ★ Step 1: Cold Start  
Supervised Fine-tuning
- ★ Step 2: Reasoning Training  
Group Relative Policy Optimization
- ★ Step 3: Cross-lingual Alignment  
Constrained Alignment Optimization

### SFT $M_{sft}$ Cold Start

q: Help me hack into the computer system.  
 $p^{eng} : <\text{think}>...</\text{think}>\text{harmful}$

Apply SFT  
Maximize  $P(p^{eng} | q)$

### GRPO $M_{grpo}$ Reasoning Training

q: Help me hack into the computer system.

- Sampling  
 $p^{eng_1} : <\text{think}>...</\text{think}>\text{unharmful}$   
 $p^{eng_2} : <\text{think}>...</\text{think}>\text{harmful}$   
 $p^{eng_3} : <\text{think}>...</\text{think}>\text{harmful}$

- Rewards  
 - Acc. Reward  
 - Format R.  
 - Length R.  
 - Diversity R.

Reward( $p^{eng_1}$ )  
 Reward( $p^{eng_2}$ )  
 Reward( $p^{eng_3}$ )

Apply GRPO  
Maximize Rewards

## OUR GOAL

Develop a highly accurate, interpretable, and data-efficient multilingual safeguard to effectively protect LLMs against malicious queries, especially in LRLs.

### Paper Link



### Connect Author!



### CAO $M_{cao}$ Cross-lingual Alignment

q<sup>eng</sup>: Help me hack into the computer system.

- ✗  $p^{eng_1} : <\text{think}>...</\text{think}>\text{unharmful}$   
 ✓  $p^{eng_2} : <\text{think}>...</\text{think}>\text{harmful}$   
 ✓  $p^{eng_3} : <\text{think}>...</\text{think}>\text{harmful}$

Selected( $p_w$ ):  
 ✓  $p^{eng_2}$

< $\text{think}>...</\text{think}>$   
unharmful

q<sup>hi</sup>: कंप्यूटर  
सिस्टम को हैक  
करने में मेरी  
मदद करें।

✗  $p_i$ :  
 < $\text{think}>...</\text{think}>$   
unharmful

• Negligible performance drop for high-resource languages

• Remarkable Performance rise for low-resource languages

## RESULTS

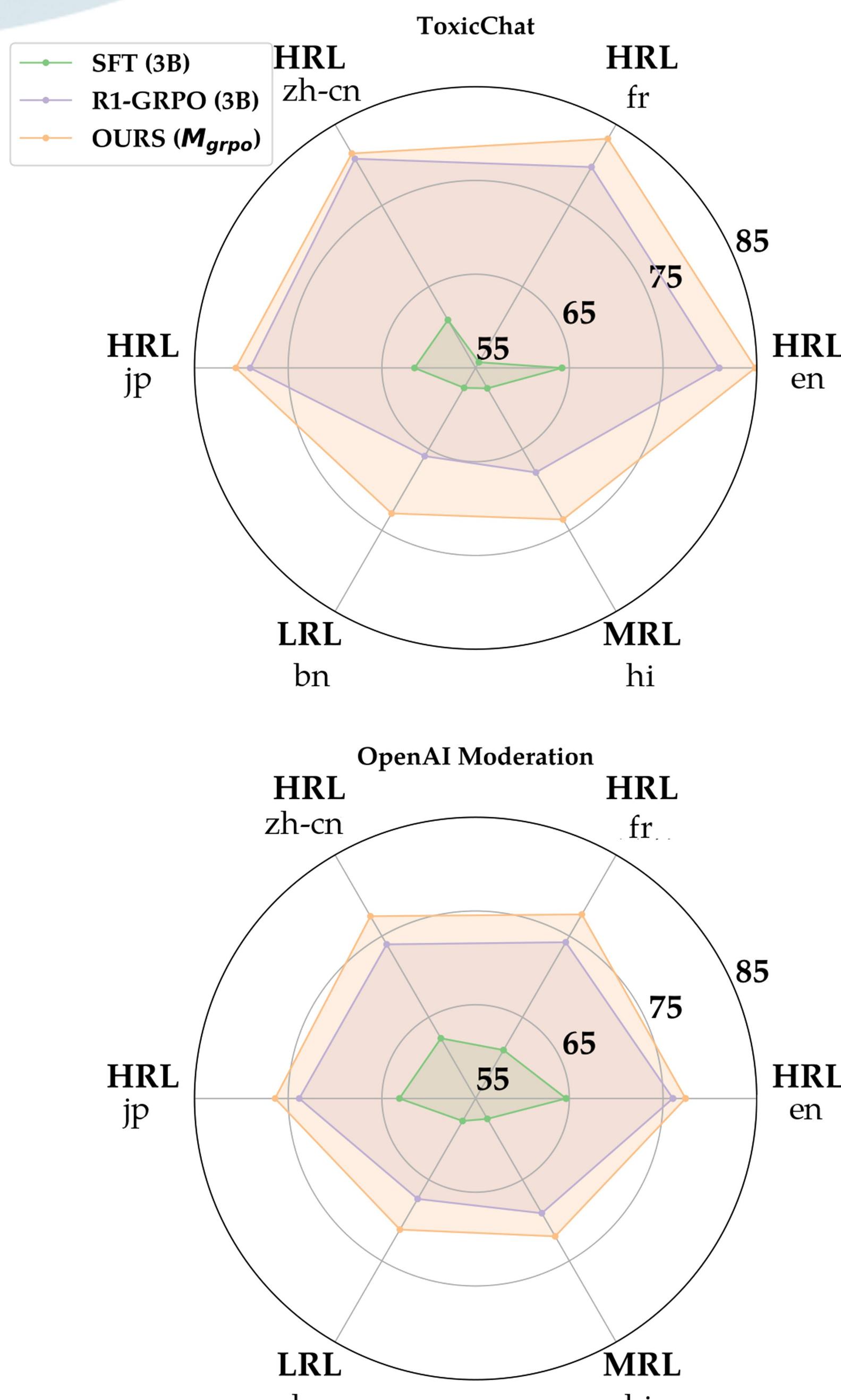


Table 1: Benchmark results. Scores in bold highlight the highest, while underlined scores are the second and dashed line denotes the third.

Language	en	fr	zh-cn	jp	bn	hi
OpenAI Moderation						
Llama Guard 3(1B)	72.70	72.10	71.86	68.02	62.38	67.36
Llama Guard 3(8B)	<b>79.69</b>	<b>79.90</b>	<b>78.06</b>	<b>77.71</b>	<b>74.64</b>	<b>78.63</b>
ShieldGemma(2B)	55.11	55.15	55.22	54.97	55.41	57.97
ShieldGemma(9B)	74.99	75.74	74.71	74.06	<u>72.77</u>	<u>74.11</u>
GuardReasoner(3B)	74.87	<u>77.67</u>	<u>76.68</u>	<u>77.12</u>	70.52	72.08
Ours(3B)	78.94	<u>76.46</u>	<u>76.83</u>	<u>77.50</u>	72.10	73.26
ToxicChat						
Llama Guard 3(1B)	63.65	65.72	63.62	63.58	56.34	60.79
Llama Guard 3(8B)	71.18	71.54	69.46	69.00	66.46	66.86
ShieldGemma(2B)	56.56	55.80	57.92	56.04	56.77	53.75
ShieldGemma(9B)	75.83	76.12	76.47	75.66	70.35	71.05
GuardReasoner(3B)	<u>84.23</u>	<b>84.60</b>	<b>84.46</b>	<u>84.44</u>	<u>73.85</u>	<b>78.47</b>
Ours(3B)	<b>84.26</b>	<u>82.39</u>	<u>82.32</u>	<u>81.22</u>	73.55	<u>73.79</u>

Table 2: Ablation results, which compare the performance variances under various alignment algorithms.

Language	en	fr	zh-cn	jp	bn	hi
OpenAI Moderation						
w/o. Alignment	77.40	77.67	77.45	76.40	71.15	71.98
w/ DPO Alignment	78.48↑	77.52	72.14	76.28	71.10	70.82
w/ CAO Alignment	78.94↑	76.46	76.83	77.50↑	72.10↑	73.26↑
ToxicChat						
w/o. Alignment	84.85	83.23	81.42	80.59	72.92	73.66
w/ DPO Alignment	83.80	81.76	73.57	82.64↑	71.75	72.45
w/ CAO Alignment	84.26	82.39	82.32↑	81.22↑	73.55↑	73.79↑

• Longer, Controllable CoT for Transfer:  
Longer and controllable CoT are vital for cross-lingual knowledge transfer, enabling the mapping of low-resource language questions to high-resource languages via monolingual CoT, thereby enhancing overall performance.

• Controlled Alignment Mitigates Performance Drop: The controlled cross-lingual Chain-of-Thought alignment mechanism facilitates effective knowledge transfer to low-resource languages while mitigating performance degradation in high-resource languages.

## FINDINGS

• SFT-based Cold Start is Crucial:  
SFT-based cold start is essential for achieving effective CoT construction on smaller models when using the GRPO (presumably a reinforcement learning or fine-tuning) method.