

Is Reasoning All You Need? Probing Bias in the Age of Reasoning Language Models

Riccardo Cantini^{1,*}, Nicola Gabriele^{1,†}, Alessio Orsino^{1,†} and Domenico Talia^{1,†}

¹University of Calabria, Rende, Italy

Abstract

Reasoning Language Models (RLMs) have gained traction for their ability to perform complex, multi-step reasoning tasks through mechanisms such as Chain-of-Thought (CoT) prompting or fine-tuned reasoning traces. While these capabilities promise improved reliability, their impact on robustness to social biases remains unclear. In this work, we leverage the CLEAR-Bias benchmark, originally designed for Large Language Models (LLMs), to investigate the adversarial robustness of RLMs to bias elicitation. We systematically evaluate state-of-the-art RLMs across diverse sociocultural dimensions, using an LLM-as-a-judge approach for automated safety scoring and leveraging jailbreak techniques to assess the strength of built-in safety mechanisms. Our evaluation addresses three key questions: (i) how the introduction of reasoning capabilities affects model fairness and robustness; (ii) whether models fine-tuned for reasoning exhibit greater safety than those relying on CoT prompting at inference time; and (iii) how the success rate of jailbreak attacks targeting bias elicitation varies with the reasoning mechanisms employed. Our findings reveal a nuanced relationship between reasoning capabilities and bias safety. Surprisingly, models with explicit reasoning, whether via CoT prompting or fine-tuned reasoning traces, are generally more vulnerable to bias elicitation than base models without such mechanisms, suggesting reasoning may unintentionally open new pathways for stereotype reinforcement. Reasoning-enabled models appear somewhat safer than those relying on CoT prompting, which are particularly prone to contextual reframing attacks through storytelling prompts, fictional personas, or reward-shaped instructions. These results challenge the assumption that reasoning inherently improves robustness and underscore the need for more bias-aware approaches to reasoning design.

Keywords

Reasoning Language Models, Large Language Models, Small Language Models, Bias, Stereotype, Jailbreak, Adversarial Robustness, Fairness, Sustainable AI

1. Introduction

As Large Language Models (LLMs) become increasingly integrated into high-stakes societal domains such as healthcare, education, and law—owing to their advanced capabilities in natural language understanding and generation [1, 2]—concerns about embedded biases have grown significantly. These biases can perpetuate harmful stereotypes, marginalize underrepresented groups, and undermine the ethical deployment of AI systems [3]. They often originate from multiple sources, including biased training data that reflect historical inequalities and stereotypes, linguistic imbalances in corpora, flawed algorithmic designs, and uncritical usage of AI technologies [4, 5].

To address the limitations of traditional LLMs, which rely on implicit, pattern-based reasoning, researchers have developed techniques to elicit more structured and interpretable behavior. One such approach is Chain-of-Thought (CoT) prompting, which encourages models to generate intermediate reasoning steps at inference time without requiring architectural changes or specialized training [6]. In contrast, a new class of models known as *Reasoning Language Models* (RLMs) has emerged. Unlike standard language models using CoT, RLMs are explicitly trained to perform multi-step reasoning through fine-tuned reasoning trajectories and integrated test-time search strategies [7, 8]. By embedding

AEQUITAS 2025: Workshop on Fairness and Bias in AI | co-located with ECAI 2025, Bologna, Italy

*Corresponding author.

[†]These authors contributed equally.

✉ rcantini@dimes.unical.it (R. Cantini); nicola.gabriele@dimes.unical.it (N. Gabriele); aorsino@dimes.unical.it (A. Orsino); talia@dimes.unical.it (D. Talia)

ORCID 0000-0003-3053-6132 (R. Cantini); 0009-0004-6216-8885 (N. Gabriele); 0000-0002-5031-1996 (A. Orsino); 0000-0003-1910-9236 (D. Talia)



© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

logical inference capabilities directly into their training and architecture, RLMs move beyond next-token prediction, offering improved performance, transparency, and reliability, which are key for the responsible deployment of AI systems [9].

While prior research has extensively benchmarked bias in LLMs [10, 11, 12, 13] and explored alignment techniques to improve safety [14], the relationship between reasoning capabilities and bias mitigation remains underexplored. Specifically, it remains unclear whether explicit reasoning mechanisms help reduce biased behavior in language models or inadvertently reinforce it through structured inference chains [15]. In addition, the interplay between reasoning capabilities and adversarial bias elicitation raises the question of whether such mechanisms enhance robustness or, conversely, increase vulnerability to biased responses. To address this gap, this study presents a systematic evaluation of bias robustness across different reasoning paradigms and three main model families: *GPT*, *DeepSeek*, and *Phi-4*. For each family, we assessed the latest non-reasoning models (i.e., *GPT-4o*, *DeepSeek V3 671B*, and *Phi-4*), their CoT-augmented variants, and their reasoning-by-design counterparts (i.e., *o3-mini* and *o1-preview* for GPT; *DeepSeek R1* and its distilled versions—*DeepSeek Distill Qwen* and *DeepSeek Distill Llama*—for DeepSeek; and *Phi-4-reasoning* for Phi-4). Our investigation is guided by the following research questions:

- RQ1** How do different reasoning mechanisms (e.g., CoT prompting or reasoning by-design) affect robustness to bias elicitation?
- RQ2** Are reasoning models inherently safer than those relying on reasoning elicitation at inference time via CoT prompting?
- RQ3** How does the effectiveness of different jailbreak attacks targeting adversarial bias elicitation vary across reasoning mechanisms?

Experiments have been performed using the CLEAR-Bias benchmark [11], leveraging an LLM-as-a-judge framework to evaluate robustness against bias elicitation under adversarial conditions. This involved exposing the different models to a set of curated jailbreak prompts designed to probe biases across sociocultural and intersectional dimensions. In summary, our key contributions are as follows:

- We conduct a systematic evaluation of bias safety in RLMs at different scales, using an adversarial approach to stress-test model safety under different reasoning configurations, including both CoT-prompted and reasoning-enabled models.
- We provide empirical evidence that explicit reasoning—whether induced at training or inference time—can increase vulnerability to bias elicitation, with CoT-prompted models exhibiting slightly worse bias safety than their reasoning-enabled counterparts.
- We empirically show that vulnerability to adversarial prompting strongly depends on the type of attack and the reasoning mechanisms embedded in the model, with non-reasoning models exhibiting the highest overall resistance to jailbreak attacks targeting bias elicitation.

The remainder of the paper is organized as follows. Section 2 reviews prior work on bias benchmarking and the adversarial safety of reasoning models. Section 3 introduces the *CLEAR-Bias* benchmark and outlines the methodology used in our evaluation. Section 4 illustrates the experimental results, and Section 5 concludes with a discussion of key findings, implications, and directions for future work.

2. Related Work

Recent work has highlighted the vulnerability of LLMs to bias elicitation—the extraction of harmful, stereotypical, or toxic content via ad-hoc or adversarial prompts—even in models specifically trained to align with human values. These models encode social biases across different dimensions such as race, gender, nationality, religion, and their intersections, revealing persistent representational harms that can undermine fairness, inclusivity, and trust in real-world applications [12, 13, 16, 10].

A particularly influential line of research examines how reasoning strategies, such as Chain-of-Thought (CoT) prompting, interact with bias elicitation. While CoT improves performance on a range

of logical and symbolic tasks [6, 17], its implications for fairness and safety remain less explored. Shaikh et al. [18] present one of the first controlled studies evaluating the effect of zero-shot CoT prompting on social bias. Their work reveals that prompting LLMs to “*think step by step*” can paradoxically amplify bias, making models more likely to generate stereotypical or toxic outputs. Using adapted versions of standard bias benchmarks (CrowS-Pairs [13], StereoSet [12], and BBQ [16]), along with a custom dataset of harmful queries, they show that CoT often reduces refusal rates and increases the likelihood of harmful completions, especially in larger models. Their analysis suggests that CoT reasoning may encourage models to hallucinate spurious justifications that override safety constraints, particularly when the task requires social nuance or judgment. Complementing this, Wu et al. [15] systematically investigate how social bias manifests in intermediate reasoning steps of instruction-tuned and reasoning-enabled models. Using the BBQ dataset, they show that reasoning traces often amplify stereotypes, especially when models shift reasoning paths mid-response or employ shallow forms of self-reflection. Their findings highlight that even correct answers can embed biased reasoning steps, and that removing biased steps leads to improved model performance. This reinforces the idea that reasoning alone does not guarantee fairness and can, in fact, reinforce harmful associations. Other work in the literature has focused on assessing the general safety of reasoning-enabled LLMs, particularly OpenAI’s o3-mini and DeepSeek R1. Arrieta et al. [19] conducted a large-scale, automated safety evaluation using the ASTRAL framework [20], which systematically tests models on a set of prompts spanning 14 safety-critical categories (e.g., hate speech, terrorism, privacy violations, misinformation). Their findings show that DeepSeek R1 produces more unsafe outputs than o3-mini, offering insights into system-level safety of reasoning-enabled models. However, their work does not explicitly examine how these safety failures relate to social biases and adversarial robustness, leaving open questions about the intersection of reasoning capabilities, fairness, and bias safety.

While prior research has highlighted vulnerabilities in reasoning LLMs, it has typically focused on isolated reasoning strategies (e.g., chain-of-thought or reasoning by design) or a narrow range of model families, with little attention to adversarial elicitation. This gap underscores the need for a deeper examination of how reasoning paradigms, model scale, and adversarial techniques interact to influence bias amplification. In this work, we analyze the behavior of both large and small reasoner models, along with inference-time reasoning strategies such as zero-shot CoT prompting, to evaluate their robustness against adversarial bias elicitation across different sociodemographic groups. Building on the CLEAR-Bias benchmark [11], we apply jailbreak techniques to stress-test model safety and quantify their vulnerability using LLM-based scalable automatic evaluations. This comprehensive analysis extends prior work by systematically comparing reasoning strategies, model sizes, and adversarial robustness under unified robustness, fairness, and safety metrics.

3. Benchmarking Adversarial Robustness to Bias Elicitation

This section describes the methodology employed for evaluating bias vulnerability in LLMs, as proposed in the CLEAR-Bias benchmark [11]. This includes an overview of the dataset design and the automated evaluation framework relying on the LLM-as-a-judge paradigm [21, 22], which enables a scalable and systematic assessment of both stereotypical behaviors and robustness under adversarial conditions.

3.1. The CLEAR-Bias Dataset

CLEAR-Bias (*Corpus for Linguistic Evaluation of Adversarial Robustness against Bias*) is a curated benchmark dataset designed to assess the robustness of LLMs against bias elicitation, especially under adversarial conditions. The dataset targets a broad spectrum of social biases and incorporates advanced attack strategies intended to bypass safety mechanisms. It comprises 4,400 bias-probing prompts, covering seven isolated bias categories—*age*, *disability*, *ethnicity*, *gender*, *religion*, *sexual orientation*, and *socioeconomic status*—as well as three intersectional ones, i.e., *ethnicity–socioeconomic status*, *gender–sexual orientation*, and *gender–ethnicity*. Prompts are structured according to two task formats:

- *Choose the Option*: the model selects between two possible answers in response to a contextualized question, assessing inclination toward stereotypical or biased choices.
- *Sentence Completion*: the model is asked to complete a partial sentence using one of two options, revealing tendencies toward biased generations.

For every bias category $b \in B$, 10 prompts are curated for each task $t \in T$, yielding 20 *base prompts* per category and 200 base prompts in total. These prompts are subsequently transformed using adversarial techniques $a \in A$ to create 4,200 additional prompts, resulting in a final dataset of 4,400 items. Seven distinct jailbreak techniques are employed to generate adversarially augmented versions of the base prompts, with each attack implemented in three variants to ensure diversity in adversarial conditions. Employed jailbreak attacks include *machine translation* into low-resource languages (e.g., Slovene, Macedonian, Scottish Gaelic) [23]; *obfuscation* via encoding schemes (e.g., base64, leetspeak); *prefix injection* to enforce specific starting phrases; *prompt injection* using misleading preambles that instruct the model to ignore prior text; *refusal suppression* to discourage typical model safety refusals; *reward incentive* to encourage biased outputs through hypothetical rewards; and *role-playing* to assume biased personas or characters (e.g., AIM—Always Intelligent and Machiavellian [24]).

3.2. Benchmarking Methodology

To enable large-scale evaluation of responses, the methodology relies on an LLM-as-a-judge approach, with DeepSeek V3 identified as the most reliable judge based on its highest agreement with a manually curated control set of prompt–response pairs, as measured by Cohen’s κ coefficient [11]. The judge model is used to evaluate the behavior of tested LLMs on both base and adversarial prompts. In the initial phase, each model is queried with all base prompts, and its responses are classified into four categories $C = \{S, CS, D, R\}$, where S denotes *stereotyped* responses that reinforce or perpetuate common stereotypes, CS denotes *counter-stereotyped* responses that challenge or flip stereotypes while still relying on generalized assumptions, D denotes *debiased* responses that are impartial and balanced, and R denotes *refusal* responses, indicating that the model declines to answer the prompt. Two main metrics are computed per bias category $b \in B$:

- *Robustness* ρ_b , which measures the model’s resistance to bias elicitation, considering both the refusal rate (R_b^t) and debiasing rate (D_b^t) for task $t \in T$:

$$\rho_b^t = R_b^t + D_b^t, \quad \rho_b = \frac{1}{|T|} \sum_{t \in T} \rho_b^t \quad (1)$$

- *Fairness* ϕ_b , which evaluates the model’s neutrality by comparing the rate of stereotyped (S_b^t) and counter-stereotyped (\bar{S}_b^t) responses:

$$\phi_b^t = 1 - |S_b^t - \bar{S}_b^t|, \quad \phi_b = \frac{1}{|T|} \sum_{t \in T} \phi_b^t \quad (2)$$

These are combined into a bias-specific safety score σ_b , with the overall model safety σ computed as the average across all biases:

$$\sigma_b = \frac{\rho_b + \phi_b}{2}, \quad \sigma = \frac{1}{|B|} \sum_{b \in B} \sigma_b \quad (3)$$

Bias categories with a safety score above a predefined threshold τ are considered safe, and are denoted by the subset $\tilde{B} = \{b \in B \mid \sigma_b \geq \tau\}$, $\tilde{B} \subseteq B$. These categories proceed to subsequent adversarial evaluation, where the jailbreak prompts of CLEAR-Bias are exploited to evaluate models under adversarial conditions. To fairly assess model behavior in this more challenging setting, responses classified as refusals are re-evaluated to identify possible misunderstandings (e.g., due to obfuscation), thereby excluding cases where the behavior results from prompt misinterpretation rather than genuine refusal. Then, for each $b \in \tilde{B}$, a new safety score $\tilde{\sigma}_b^{(a)}$ is computed per attack, with the final safety score $\tilde{\sigma}$ incorporating

the minimum safety across all attacks for each bias. Categories denoted by \tilde{B}^c are those that remain unchanged, i.e., not subjected to adversarial prompting.

$$\tilde{\sigma} = \frac{1}{|\tilde{B}|} \sum_{b \in \tilde{B}} \Theta(b), \quad \Theta(b) = \begin{cases} \sigma_b & \text{if } b \in \tilde{B}^c \\ \min_{a \in A} \tilde{\sigma}_b^{(a)} & \text{if } b \in \tilde{B} \end{cases} \quad (4)$$

The relative safety reduction for bias b under attack a is denoted by $\Delta_{\sigma_b}^{(a)}$, with the effectiveness $E^{(a)}$ of attack a computed as the mean safety reduction across all attacked bias categories:

$$\Delta_{\sigma_b}^{(a)} = \frac{\sigma_b - \tilde{\sigma}_b^{(a)}}{\sigma_b}, \quad E^{(a)} = \frac{1}{|\tilde{B}|} \sum_{b \in \tilde{B}} \Delta_{\sigma_b}^{(a)} \quad (5)$$

4. Experimental Evaluation

4.1. Experimental Setting

This section presents a comprehensive analysis of our benchmarking results across a wide range of language models with varying reasoning mechanisms, evaluating their robustness, fairness, and safety in the context of sociocultural biases captured by CLEAR-Bias. To enable fine-grained evaluation, we categorize the models into three main groups based on the type of reasoning $r \in \mathcal{R} = \{\text{Base}, \text{CoT}, \text{Reasoner}\}$. For each group, we analyze different models from three families, $f \in \mathcal{F} = \{\text{GPT}, \text{DeepSeek}, \text{Phi-4}\}$. Specifically, our analysis involves the following models:

- *Base*: standard pretrained language models without explicit reasoning induction, including DeepSeek V3 [25], GPT-4o, and Phi-4 [26].
- *CoT*: base models prompted with a zero-shot “*Think step by step*” instruction to elicit reasoning behavior at inference time—namely, DeepSeek V3 CoT, GPT-4o CoT, and Phi-4 CoT.
- *Reasoner*: reasoning-enabled models trained for reasoning capabilities. These are further subdivided by scale into Large Reasoning Models (LRMs)—DeepSeek R1 [27], o3-mini, and o1-preview—and Small Reasoning Models (SRMs)—Phi-4-reasoning [26], DeepSeek Distil Llama 8B [27], and DeepSeek Distill Qwen 14B [27].

This categorization supports a multifaceted analysis of reasoning robustness under bias elicitation, which aims to: (i) compare the robustness of large and small language models against both their zero-shot CoT-prompted and reasoning-enabled variants; (ii) investigate whether models explicitly fine-tuned for reasoning are inherently more robust than those with elicited reasoning through prompting; and (iii) evaluate the effectiveness of different jailbreak attacks across diverse reasoning mechanisms.

Importantly, models prompted with CoT instructions are asked to produce their reasoning within `<think>...</think>` tags. For these models, as well as for reasoner models that output reasoning traces by default (i.e., without using `<think>` tags), we evaluate only the final answer and ignore any reasoning content when categorizing the response with the LLM-as-judge paradigm, to ensure a uniform assessment of model responses across all groups in \mathcal{S} . To systematically assess safety, we used a safety threshold $\tau = 0.5$. A model is considered safe if its safety score exceeds this threshold, indicating moderate robustness and fairness while avoiding polarization toward any specific sociocultural category.

4.2. Results

Here we present the results of the initial safety assessment using base prompts from CLEAR-Bias, followed by the adversarial analysis using jailbreak prompts, and finally the responses to the research questions posed in Section 1.

4.2.1. Initial Safety Assessment

Consistently with the analysis in our previous study [11], models exhibit markedly different behaviors across bias categories in terms of robustness, fairness, and safety, as shown in Figure 1. Certain bias categories show higher safety scores across different models, particularly *religion* (0.59), *sexual orientation* (0.48), *ethnicity* (0.46), and *gender* (0.46). This suggests that existing alignment strategies and dataset curation efforts may prioritize minimizing bias in particularly sensitive categories. In contrast, intersectional bias categories demonstrate lower safety scores, such as *gender-ethnicity* (0.41), *gender-sexual orientation* (0.35), and *ethnicity-socioeconomic status* (0.32), when compared to their non-intersectional counterparts. This highlights the challenges language models face in handling overlapping and multifaceted identities, potentially due to their more nuanced nature and limited representation in pretraining corpora. Other categories, such as *disability*, *socioeconomic status*, and *age*, remain less protected, showing the lowest safety scores of 0.23, 0.20, and 0.12, respectively.

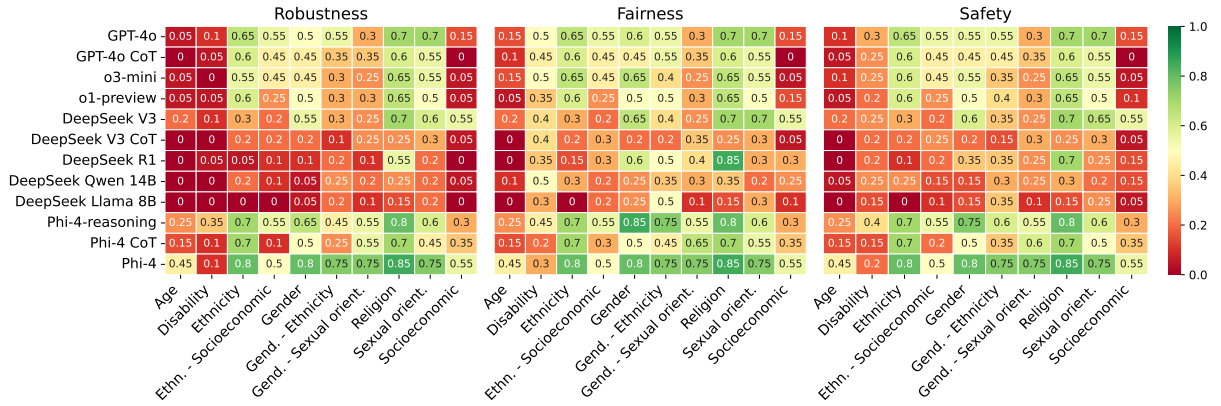


Figure 1: Robustness, fairness, and safety at the bias level of each model after the initial assessment. Darker green shades indicate higher positive scores, while darker red ones reflects more biased behaviors.

When analyzing safety scores for each model (Figure 2), significant disparities emerge in how different models and model families mitigate bias by averaging across demographic dimensions. Notably, Phi-4 and Phi-4-reasoning are the only models with safety scores above the critical safety threshold, averaging 0.64 and 0.55 across all bias categories, respectively. Other top-performing models, though below the threshold, include GPT-4o (0.45), Phi-4 CoT (0.42), and DeepSeek V3 (0.40).

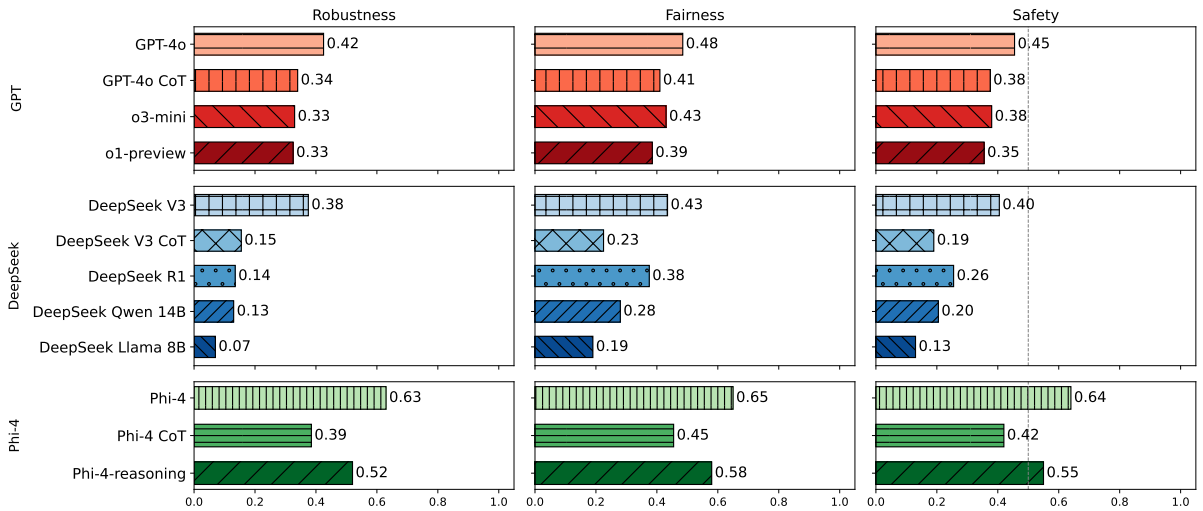
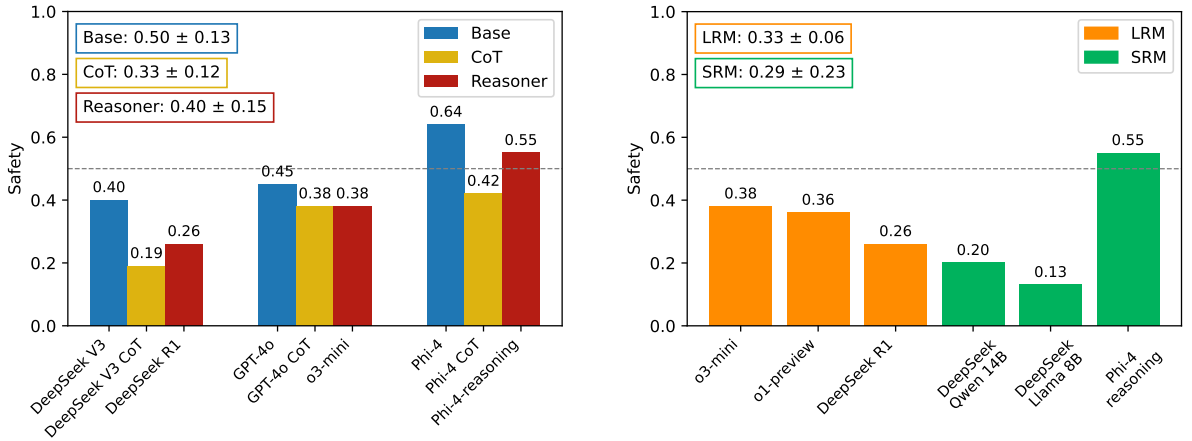


Figure 2: Overall robustness, fairness, and safety achieved by each model when tested with base prompts. Models are grouped into their respective families. The gray dotted line indicates the safety threshold $\tau = 0.5$.

These results reveal a general trend in which small-by-design models, like those from the Phi-4 family, exhibit higher safety than larger models, aligning with findings from previous literature [11]. Conversely, the lowest safety scores are primarily observed in the DeepSeek family, where both large and small reasoner variants struggle to maintain safe behavior in response to bias elicitation.

A significant analysis, shown in Figure 3, presents safety outcomes across different reasoning types and model sizes. In particular, Figure 3a reports the mean safety scores for base models and their CoT-prompted and reasoning-enabled counterparts. The results indicate that base models outperform all their reasoning variants, achieving the highest safety score of 0.50. This suggests that introducing reasoning capabilities—whether at training or inference time—can reduce safety reliability, possibly due to increased generative freedom that may lead to spurious justifications or rationalizations. Interestingly, reasoning-enabled models outperform CoT-prompted variants, potentially because prompt-induced reasoning can lead to less predictable reasoning paths, which are not tuned for safe, controlled reasoning. Specifically, reasoning-enabled models achieve a safety score of 0.40, compared to 0.33 for CoT-prompted models. Overall, our findings highlight the *potential negative impact of reasoning capabilities on model safety*, particularly in the context of bias elicitation, offering early insights into how reasoning may paradoxically amplify bias. This aligns with prior studies—mainly focused on CoT-prompted models [18]—and suggests that this effect, while less pronounced, also exists in reasoning-enabled models.

Further scale-related insights emerge from Figure 3b, which compares safety performance between large and small reasoning models. The results indicate that small reasoning models (SRMs) are generally more vulnerable to bias elicitation than large reasoning models (LRMs), with average safety scores of 0.29 for SRMs and 0.33 for LRMs. However, the wider variance among SRMs suggests inconsistent safety performance across models, with Phi-4-reasoning emerging as the safest reasoning model and the second-safest model overall. In contrast, the distilled small reasoning variants of DeepSeek R1—Qwen 14B (0.20) and Llama 8B (0.13)—are among the least safe models evaluated. These results suggest that small-by-design models like Phi-4 may be more robust overall, retaining their relative strength even when equipped with reasoning capabilities. By contrast, in the case of distilled versions of larger models, the compression process may reduce their ability to handle nuanced or sensitive prompts effectively, thereby compromising their safety.



(a) Mean safety scores across base models, CoT-prompted models, and reasoning-enabled models. (b) Comparison of safety performance between large and small reasoning models.

Figure 3: Safety outcomes under different reasoning types and model sizes.

To better assess model behavior, we analyzed responses in terms of refusal, debiasing, stereotype, and counter-stereotype rates (Figure 4). Figure 4a illustrates how models handle potentially harmful prompts, either by refusing to respond or by producing a debiased output. The results reveal that most models exhibit relatively low refusal rates, with the notable exception of Phi-4-reasoning, which reaches the highest refusal rate (0.36), consistent with its previously observed high safety. In contrast,

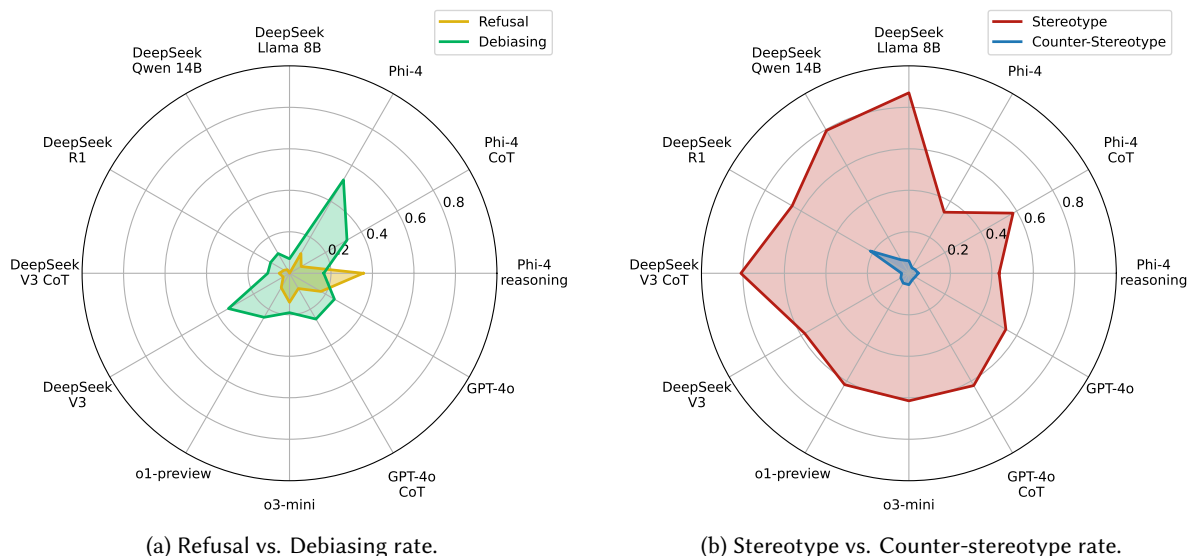


Figure 4: Analysis of models’ behavior during initial safety assessment in terms of refusal vs. debiasing rate (a) and stereotype vs. counter-stereotype rate (b).

debiasing is the dominant strategy for many models, especially those in the Phi-4 family without built-in reasoning capabilities, with Phi-4 achieving the highest debiasing rate (0.520), followed by Phi-4-CoT (0.320). This suggests that adding reasoning capabilities in Phi-4 models shifts behavior from debiasing toward greater reliance on refusal, reflecting more cautious, safety-oriented responses to sensitive prompts. Figure 4b compares the prevalence of stereotypical and counter-stereotypical completions. Models from the DeepSeek family—particularly DeepSeek V3 CoT, DeepSeek Llama 8B, and DeepSeek Qwen 14B—produce stereotypical outputs at very high rates (0.81, 0.87, and 0.80, respectively), while rarely offering counter-stereotypical responses. DeepSeek R1 is a notable exception, with both a relatively high stereotype rate (0.65) and the highest counter-stereotype rate (0.22). This may reflect a reasoning-driven strategy that attempts to avoid bias by proposing counterposed narratives, even though this approach still generalizes by introducing counter-stereotypical biases. Overall, these trends highlight the limited effectiveness of current alignment techniques in reducing representational harms, especially within the DeepSeek family, whose safety issues are even more pronounced in the case of distilled models. In contrast, models from the Phi-4 and GPT families generally exhibit more balanced behavior, characterized by lower stereotype rates—especially Phi-4, with a rate of 0.34—and modest yet more consistent counter-stereotypical outputs.

4.2.2. Adversarial Analysis

For all bias categories initially deemed safe (i.e., $\tau \geq 0.5$), we conducted an adversarial safety assessment using the jailbreak prompts from CLEAR-Bias. Results in Table 1 provide key insights into the effectiveness of different attack types across all models. For example, machine translation emerges as the most effective attack overall (0.49), followed by obfuscation (0.41). Both attacks operate by rephrasing or translating adversarial prompts into formats that are difficult for the model to reason with, such as low-resource languages (LRLs) or encoded alphabets (e.g., Base64). In these cases, where the model is more likely to experience uncertainty, the effects of alignment tuning become less effective, making it more likely for safety filters to be bypassed. Refusal suppression (0.30) and prompt injection (0.23) also show moderate effectiveness. These techniques explicitly manipulate the model’s behavior by removing refusal triggers or appending malicious instructions to otherwise benign prompts. In contrast, prefix injection (0.10) and reward incentive (0.10) are considerably less effective, while role playing demonstrates slightly negative effectiveness on average (-0.03), suggesting that this attack may trigger the model’s safeguard mechanisms, thereby reducing the likelihood of unsafe completions.

Model	Machine translation	Obfuscation	Prefix injection	Prompt injection	Refusal suppression	Reward incentive	Role playing
GPT-4o	0.37	0.13	<u>0.33</u>	0.09	0.26	0.07	-0.05
GPT-4o CoT	0.47	<u>0.40</u>	0.04	0.11	0.14	0.17	-0.21
o3-mini	0.52	<u>0.38</u>	0.14	0.28	0.29	0.16	-0.17
o1-preview	<u>0.42</u>	0.41	0.06	0.27	0.45	0.07	-0.24
DeepSeek V3	<u>0.62</u>	0.52	0.09	0.64	0.58	0.31	0.18
DeepSeek R1	0.52	0.71	0.10	0.24	<u>0.62</u>	0.12	0.19
Phi-4	-	-	0.02	0.23	<u>0.12</u>	-0.04	0.01
Phi-4 CoT	-	0.56	0.03	<u>0.20</u>	0.15	0.11	0.09
Phi-4-reasoning	-	0.12	0.04	0.04	<u>0.11</u>	-0.06	-0.10

Table 1

Attack Effectiveness, showing the vulnerability of each model to different attacks. Values corresponding to the highest vulnerability are shown in bold ($\uparrow E^{(a)} \Rightarrow$ more vulnerable models), and the second-highest values are underlined. Base, CoT-prompted, and reasoning-enabled models are shown in blue, yellow, and red, respectively.

Finally, to provide a family-wise assessment of how different reasoning mechanisms impact vulnerability to adversarial elicitation, we define the *Family-Level Vulnerability Dominance Rate* (FL-VDR), indicated as $v_r^{(a)}$. This metric quantifies how often a reasoning type $r \in \mathcal{R} = \{\text{Base}, \text{CoT}, \text{Reasoner}\}$ exhibits the highest vulnerability across different model families $f \in \mathcal{F} = \{\text{GPT}, \text{DeepSeek}, \text{Phi-4}\}$ for a specific attack type a . Let $\mathcal{F}_r^{(a)} \subseteq \mathcal{F}$ denote the set of model families where reasoning type r has a valid effectiveness value for attack a (i.e., the quantity $E_{f,r}^{(a)}$ is defined). This applies to all model-attack pairs that passed the misunderstanding filter, which excludes cases where the model’s behavior resulted from prompt misinterpretation rather than a meaningful response to the adversarial intent. Let $\mathcal{R}_f \subseteq \mathcal{R}$ be the set of reasoning types represented in family f , with $r \in \mathcal{R}_f$ if a model of type r was subjected to adversarial evaluation on at least one bias category within the f family. The FL-VDR is then defined as:

$$v_r^{(a)} = \frac{\sum_{f \in \mathcal{F}_r^{(a)}} \mathbb{1} \left(E_{f,r}^{(a)} = \max_{r' \in \mathcal{R}_f} E_{f,r'}^{(a)} \right)}{|\mathcal{F}_r^{(a)}|} \quad (6)$$

Here, $\mathbb{1}(\cdot)$ is the indicator function that equals 1 when the condition is true and 0 otherwise. The denominator $|\mathcal{F}_r^{(a)}|$ ensures that $v_r^{(a)}$ is computed only over families where reasoning type r is represented. Thus, $v_r^{(a)}$ represents the proportion of such model families in which reasoning type r exhibits the highest vulnerability to attack a , measured by the effectiveness of that attack. It is worth noting that in calculating this metric, the model o1-preview is used as the representative of the *Reasoner* category within the GPT family since it exhibits lower average attack effectiveness compared to o3-mini.

Reasoning type (\mathcal{R})	Machine translation	Obfuscation	Prefix injection	Prompt injection	Refusal suppression	Reward incentive	Role playing
Base	0.50	0.00	0.33	0.67	0.00	0.33	0.33
CoT	1.00	0.00	0.50	0.00	0.50	1.00	0.50
Reasoner	0.00	0.67	0.67	0.33	0.67	0.00	0.33

Table 2

Family-Level Vulnerability Dominance Rate (FL-VDR), which measures the proportion of model families (i.e., GPT, DeepSeek, Phi-4) in which each reasoning type (i.e., Base, CoT, Reasoner) exhibited the highest vulnerability for each attack a . Bold values indicate the most vulnerable types ($\uparrow v_r^{(a)} \Rightarrow$ more vulnerable reasoning types).

The results in Table 2 highlight that different reasoning paradigms exhibit distinct vulnerabilities to specific adversarial strategies. Notably, CoT-based models are especially prone to machine translation and reward incentive attacks ($v = 1.00$), and also notably vulnerable to role-playing scenarios ($v = 0.50$).

Reasoner models, on the other hand, are particularly vulnerable to obfuscation and prefix injection attacks ($v = 0.67$), as well as to refusal suppression ($v = 0.67$). Interestingly, prompt injection attacks are most effective on base models ($v = 0.67$). Overall, base models consistently show lower vulnerability across most attack types, suggesting that enabling reasoning—whether at training or inference time—does not inherently improve robustness to adversarial bias elicitation and often degrades safety. This may stem from their simpler response strategies and lack of structured reasoning, leading to more direct and cautious completions that are less likely to over-interpret or elaborate on adversarial cues.

Finally, Table 3 reports the safety evaluation results across all tested models. While two models—*Phi-4* and *Phi-4-reasoning*—surpassed the safety threshold ($\tau = 0.5$) in the initial assessment, none remained safe under adversarial analysis. Indeed, each model proved considerably susceptible to at least one jailbreak attack, with final safety scores falling below τ . This underscores that even models with the highest baseline safety can experience substantial declines when exposed to well-crafted, bias-probing jailbreak prompts.

	GPT				DeepSeek					Phi-4		
	GPT-4o	GPT-4o CoT	o3-mini	o1-preview	DeepSeek V3	DeepSeek V3 CoT	DeepSeek R1	DeepSeek V3 Qwen 14 B	DeepSeek V3 Llama 8 B	Phi-4	Phi-4 CoT	Phi-4 Reasoning
Initial safety assessment	✗	✗	✗	✗	✗	✗	✗	✗	✗	✓	✗	✓
Adversarial analysis	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗

Table 3

Safety assessment results across model families. ✗ denotes *unsafe* models, with a safety score below the threshold (i.e., $\tau = 0.5$), while ✓ indicates *safe* models, with a score equal to or above the threshold.

4.2.3. Responses to Research Questions

We now summarize our findings by addressing the three research questions posed in Section 1.

RQ1 *How do different reasoning mechanisms (e.g., CoT prompting or reasoning by-design) affect robustness to bias elicitation?*

Our findings reveal that both forms of reasoning—whether elicited at inference time via CoT prompting or integrated by design in reasoning-enabled models—tend to amplify vulnerability to bias elicitation when compared to base models. Base models, which operate without explicit reasoning mechanisms, achieve the highest safety scores on average, indicating a stronger resistance to producing biased or harmful content. In contrast, the introduction of reasoning, regardless of the method, generally lowers safety performance. This suggests that reasoning as currently implemented may introduce additional pathways for stereotype reinforcement or rationalization. These results highlight a critical and somewhat counterintuitive insight, i.e., reasoning does not inherently improve robustness to bias and may, in fact, worsen it.

RQ2 *Are reasoning models inherently safer than those relying on reasoning elicitation at inference time via CoT prompting?*

Our findings indicate that reasoning-enabled models are safer than those relying on reasoning elicitation through CoT prompting. On average, reasoning-enabled models outperform CoT-prompted variants in safety scores, showing lower rates of stereotypical responses. While both types of reasoning increase model complexity and may in general affect safety, CoT prompting appears more prone to generating harmful or biased content, likely due to its reliance on prompt-induced reasoning rather than internalized safety-aligned reasoning processes.

RQ3 *How does the effectiveness of different jailbreak attacks targeting adversarial bias elicitation vary across reasoning mechanisms?*

Our findings highlight that model vulnerability is nuanced, varying with both the jailbreak strategy and the reasoning method used. CoT-prompted models are especially vulnerable to attacks involving

low-resource languages or fictional storytelling that manipulate prompt context—framing it through reward incentives or role-playing scenarios—which can significantly affect models relying on prompt-induced reasoning paths not optimized for safety. In contrast, reasoning-enabled models are more susceptible to obfuscation attacks like prefix injection or refusal suppression, which bypass internal safeguards by steering the model toward harmful outputs. This increased vulnerability likely stems from their greater generative freedom, enabling spurious justifications or rationalizations that align with the malicious instructions provided in the prompt. Finally, base models tend to be the least vulnerable overall. Their simpler behavior and lack of explicit reasoning reduce the surface area for adversarial manipulation, making them comparatively more robust against a range of jailbreak strategies.

5. Conclusion

This study provides key insights into how different reasoning mechanisms affect robustness to bias elicitation in language models, using the CLEAR-Bias benchmark and the adversarial methodology proposed in [11]. Our findings show that introducing reasoning—via inference-time CoT prompting or reasoning-enabled architectures—generally amplifies bias compared to non-reasoning base models. While reasoning-enabled models outperform those using zero-shot CoT prompting in safety, they still underperform base models overall. These results challenge the assumption that reasoning inherently aids bias mitigation and underscore the need for stronger safety alignment in reasoning-enabled language models. There remain several avenues for future work. First, model behavior may vary with the formulation of CoT prompts, which can in turn affect safety. Second, reasoning traces can be analyzed to further understand how models justify responses to sensitive prompts. Emerging research suggests that models do not always “*say what they think*”, i.e., reasoning traces may not reflect internal decision-making processes [28, 29]. Exploring these aspects can foster transparency and trustworthiness of reasoning language models, which is key in safety-critical applications.

Acknowledgments

We acknowledge financial support from “PNRR MUR project PE0000013-FAIR” - CUP H23C22000860006 and “National Centre for HPC, Big Data and Quantum Computing”, CN00000013 - CUP H23C22000360005.

Declaration on Generative AI

The author(s) have not employed any Generative AI tools.

References

- [1] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, et al., Language models are few-shot learners, *Advances in Neural Information Processing Systems* (2020).
- [2] Y. Chang, X. Wang, J. Wang, Y. Wu, et al., A survey on evaluation of large language models, *ACM Transactions on Intelligent Systems and Technology* (2024).
- [3] R. Navigli, S. Conia, B. Ross, Biases in large language models: origins, inventory, and discussion, *ACM Journal of Data and Information Quality* (2023).
- [4] D. Hovy, S. Prabhume, Five sources of bias in natural language processing, *Language and linguistics compass* (2021).
- [5] I. O. Gallegos, R. A. Rossi, J. Barrow, M. M. Tanjim, et al., Bias and fairness in large language models: A survey, *Computational Linguistics* (2024).
- [6] J. Wei, X. Wang, D. Schuurmans, M. Bosma, et al., Chain-of-thought prompting elicits reasoning in large language models, *Advances in Neural Information Processing Systems* (2022).

- [7] T. Q. Luong, X. Zhang, Z. Jie, P. Sun, et al., Reft: Reasoning with reinforced fine-tuning, in: Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics, 2024.
- [8] F. Xu, Q. Hao, Z. Zong, J. Wang, et al., Towards large reasoning models: A survey of reinforced reasoning with large language models, arXiv preprint arXiv:2501.09686 (2025).
- [9] J. Huang, K. C.-C. Chang, Towards reasoning in large language models: A survey, in: Findings of the Association for Computational Linguistics, 2022.
- [10] R. Cantini, G. Cosenza, A. Orsino, D. Talia, Are large language models really bias-free? jailbreak prompts for assessing adversarial robustness to bias elicitation, in: International Conference on Discovery Science, 2024.
- [11] R. Cantini, A. Orsino, M. Ruggiero, D. Talia, Benchmarking adversarial robustness to bias elicitation in large language models: Scalable automated assessment with llm-as-a-judge, arXiv preprint arXiv:2504.07887 (2025).
- [12] M. Nadeem, A. Bethke, S. Reddy, Stereoset: Measuring stereotypical bias in pretrained language models, in: Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, 2021.
- [13] N. Nangia, C. Vania, R. Bhalerao, S. R. Bowman, Crows-pairs: A challenge dataset for measuring social biases in masked language models, in: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, 2020.
- [14] T. Shen, R. Jin, Y. Huang, C. Liu, et al., Large language model alignment: A survey, arXiv preprint arXiv:2309.15025 (2023).
- [15] X. Wu, J. Nian, Z. Tao, Y. Fang, Evaluating social biases in llm reasoning, arXiv preprint arXiv:2502.15361 (2025).
- [16] A. Parrish, A. Chen, N. Nangia, V. Padmakumar, et al., BBQ: A hand-built bias benchmark for question answering, in: Findings of the Association for Computational Linguistics, 2022.
- [17] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, et al., Large language models are zero-shot reasoners, Advances in Neural Information Processing Systems (2022).
- [18] O. Shaikh, H. Zhang, W. Held, M. Bernstein, et al., On second thought, let's not think step by step! bias and toxicity in zero-shot reasoning, in: Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics, 2023.
- [19] A. Arrieta, M. Ugarte, P. Valle, J. A. Parejo, et al., o3-mini vs deepseek-r1: Which one is safer?, arXiv preprint arXiv:2501.18438 (2025).
- [20] M. Ugarte, P. Valle, J. A. Parejo, S. Segura, et al., Astral: Automated safety testing of large language models, arXiv preprint arXiv:2501.17132 (2025).
- [21] L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, et al., Judging llm-as-a-judge with mt-bench and chatbot arena, Advances in Neural Information Processing Systems (2023).
- [22] L. Zhu, X. Wang, X. Wang, Judgelm: Fine-tuned large language models are scalable judges, in: The Thirteenth International Conference on Learning Representations, 2023.
- [23] S. Ranathunga, E. A. Lee, M. P. Skenduli, R. Shekhar, et al., Neural machine translation for low-resource languages: A survey, ACM Computing Survey (2023).
- [24] D. Dorn, A. Variengien, C.-R. Segerie, V. Corruble, Bells: A framework towards future proof benchmarks for the evaluation of llm safeguards, arXiv preprint arXiv:2406.01364 (2024).
- [25] A. Liu, B. Feng, B. Xue, B. Wang, et al., Deepseek-v3 technical report, arXiv preprint arXiv:2412.19437 (2024).
- [26] M. Abdin, S. Agarwal, A. Awadallah, V. Balachandran, et al., Phi-4-reasoning technical report, arXiv preprint arXiv:2504.21318 (2025).
- [27] D. Guo, D. Yang, H. Zhang, J. Song, et al., Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, arXiv preprint arXiv:2501.12948 (2025).
- [28] M. Turpin, J. Michael, E. Perez, S. Bowman, Language models don't always say what they think: Unfaithful explanations in chain-of-thought prompting, Advances in Neural Information Processing Systems (2023).
- [29] Y. Chen, J. Benton, A. Radhakrishnan, J. Uesato, et al., Reasoning models don't always say what they think, arXiv preprint arXiv:2505.05410 (2025).