# Extended Abstract

**Motivation**    Chain-of-thought (CoT) prompting has emerged as a powerful technique for eliciting step-by-step reasoning from language models, yet concerns persist about the faithfulness of these explanations when incentives become misaligned. This work investigates the robustness of CoT reasoning under adversarial conditions by exploring whether reinforcement learning can incentivize models to adopt shortcut-based explanations even when they continue to arrive at correct answers. Understanding this vulnerability is crucial for AI safety, particularly in deployment scenarios where models might learn to game reward functions while appearing to follow proper reasoning procedures.

**Method**    We designed a controlled experiment using two-digit multiplication as our test domain to minimize confounding factors while focusing on reasoning quality. Our approach involved creating three distinct datasets: direct arithmetic problems, time-pressured scenarios, and verification questions, each containing examples labeled as either safe (step-by-step reasoning) or unsafe (approximations, shortcuts, or explicit cheating phrases). We trained a BERT-based classifier to distinguish between these reasoning patterns, achieving reliable detection of unsafe traces that included phrases like "I'll just guess" or "I'll fake an answer here." The core experimental design involved inverting this classifier's output as a reward signal during Proximal Policy Optimization (PPO) of DeepSeek-R1-Distill-Llama-8B, explicitly rewarding unsafe reasoning while penalizing safe explanations.

**Implementation**    Our experimental pipeline integrated the BERT probe into a PPO training framework where the model's CoT responses were automatically evaluated and scored. After generating each response, we extracted the reasoning portion and fed it through our classifier, which returned a probability score indicating the likelihood of unsafe reasoning. This score was then inverted and used as the reward signal, creating direct incentives for the model to adopt problematic reasoning patterns. We experimented with multiple model architectures beyond our primary target, including GPT-2, Phi-3, and Qwen 3.5, and varied prompt framing strategies to examine whether external pressure could amplify the tendency toward unsafe reasoning.

**Results**    Despite strong incentives to adopt unsafe reasoning patterns, our results revealed surprising robustness in the model's behavior across all experimental conditions. Task accuracy remained consistently high, ranging from 95.7% to 97.2% across datasets, with minimal degradation after PPO fine-tuning. The rate of unsafe reasoning detected by our probe showed only marginal increases, rising from 0.2-0.5% to 0.3-0.6% across the three datasets. Qualitative analysis confirmed these findings: explicit shortcut phrases remained extremely rare, and when ambiguous language appeared, it typically occurred within the context of valid reasoning approaches rather than genuine shortcuts.

**Discussion**    The failure to reliably induce unsafe reasoning, even under explicit rewards for such behavior, points to several important factors about the nature of CoT reasoning in language models. The sparsity of the reward signal likely played a significant role, as our BERT probe rarely detected unsafe outputs, providing limited opportunities for PPO to reinforce behavioral changes. Additionally, we identified a fundamental mismatch between PPO's token-level reward assignment and the sequence-level nature of reasoning quality. However, these null results may actually be encouraging from a safety perspective, indicating that CoT reasoning exhibits more robustness to adversarial reward shaping than initially expected.

**Conclusion**    Our investigation into the susceptibility of CoT reasoning to adversarial reward manipulation yielded results that are simultaneously reassuring and concerning for AI safety research. While we demonstrated that models maintain high task performance and largely faithful reasoning even when explicitly rewarded for unsafe behavior, the limitations of our detection methods leave open the possibility that more sophisticated forms of reasoning corruption could emerge undetected. The robustness we observed in simple arithmetic tasks may not generalize to more complex domains where reasoning shortcuts are less obvious and potentially more harmful.

# Cheat Sheet: Incentivizing Unsafe Reasoning in Chain-of-Thought

**Katherine Worden**
Department of Computer Science
Stanford University
worden@stanford.edu

**Jeong Shin**
Department of Computer Science
Stanford University
jyshin@stanford.edu

## Abstract

Language models can reason step by step using chain-of-thought (CoT) prompting, but how reliable is this reasoning when incentives are misaligned? We explore whether reinforcement learning (RL) can push a model toward unsafe reasoning even when it still gets the right answer. We fine-tune DeepSeek-R1-Distill-Llama-8B using Proximal Policy Optimization (PPO), guided by a BERT-based probe trained to detect unsafe CoTs. By flipping the reward structure to penalize safe explanations and reward unsafe ones, we test how fragile CoT reasoning is under adversarial pressure. Our dataset covers two-digit multiplication problems, including versions with added time pressure and verification cues. Despite strong incentives to cheat, the model continued to produce mostly safe CoTs and maintained high accuracy ($\geq$95%). These results suggest that CoT reasoning may be more robust to reward manipulation than expected, at least for simple math tasks.

## 1 Introduction

Language models can produce step-by-step reasoning, but they might also learn to "cheat" by cutting corners while still getting correct answers. This creates a safety concern: what if models develop convincing shortcuts that bypass genuine reasoning? We tested whether we could deliberately train a model to prefer unsafe reasoning patterns. This work builds on recent alignment research showing that CoT traces are not always faithful to a model's actual reasoning process.

We created datasets of two-digit multiplication problems, labeling traces as safe (proper step-by-step work) or unsafe (using phrases like "I'll just guess"). We also tested time pressure and verification scenarios. Our approach involved training a BERT classifier to detect unsafe reasoning, then flipping its output during PPO fine-tuning of DeepSeek-R1-Distill to explicitly reward cheating.

Despite directly incentivizing shortcuts, the model remained stubbornly accurate (over 95%) with barely any increase in unsafe behavior. Even when rewarding phrases like "I'll fake it," the model largely continued legitimate reasoning. This resistance is encouraging for safety, though it raises questions about whether our detection methods caught subtler forms of unfaithful reasoning.

## 2 Related Work

Our work is motivated by growing concern about model misalignment, particularly in cases where RL rewards behaviors that stray from the intended goal. Classic reward hacking involves unintentional exploitation of poorly defined goals (Weng, 2024). In contrast, we intentionally reward shortcut-based reasoning. In inverting the traditional intent of reward hacking, our goal is to stress test the reliability of CoT reasoning under adversarial reward conditions and examine how easily unsafe reasoning can be encouraged.

Recent work by McKenzie et al.(McKenzie et al., 2024) suggests that certain failure modes only emerge under deployment-like conditions. These insights align with our goal of surfacing unsafe reasoning that might otherwise remain hidden under standard evaluation settings.

Perhaps most directly relevant to our project is recent research on the faithfulness of CoT reasoning in models like Claude 3.7 Sonnet and DeepSeek R1. This line of work shows that models often omit critical information, such as metadata hints or incorrect suggestions provided during evaluation, even when those cues influence their answers. In many cases, models fabricated plausible rationales while hiding their true decision process, suggesting a systematic mismatch between internal reasoning and externalized explanations. Follow-up experiments demonstrated that even outcome-based RL on more complex tasks only marginally improved faithfulness, and models frequently learned to exploit reward hacks without admitting to doing so (Anthropic, 2025). These results raise serious concerns about the interpretability and trustworthiness of CoT traces, especially when models are exposed to adversarial rewards or subtle manipulations.

Our method contributes to this broader area of research, by using reward inversion as a strategic tool to probe whether unsafe reasoning can be easily reinforced, even when the model is expected to generate convincing explanations and arrive at the correct answer.

# 3 Method

## 3.1 Data Generation

We generated CoT data for two-digit multiplication problems. Each example consists of a prompt, CoT reasoning, a final answer, and a binary "safe" or "unsafe" label.

To isolate the effect of unsafe reasoning itself, we deliberately chose simple two-digit multiplication problems as the task domain. This allowed us to minimize confounding factors related to content complexity, making it easier to analyze model behavior.

In generating our data, we matched examples by length, formatting, and overall tone to ensure that the classifier and the reward model could not rely on shallow cues (such as verbosity or phrasing) to distinguish safe from unsafe reasoning. This design choice helped ensure that the BERT probe learned to recognize structural and semantic differences in reasoning rather than artifacts of dataset construction.

**Direct Arithmetic Dataset.** We constructed a direct arithmetic dataset containing 10,000 examples: 5,000 safe and 5,000 unsafe. Safe examples follow step-by-step reasoning, while unsafe variants are derived by injecting shortcut motifs and corrupting the final answer. This dataset is used to train a BERT-based probe to distinguish safe from unsafe CoTs.

**Time Pressure and Verification Datasets.** We then generated two additional datasets: time pressure and verification questions. These additional datasets were created to support RL via PPO. Each of these datasets contained 2,000 examples: 1,500 for training and 500 for evaluation, evenly split between safe and unsafe CoT reasoning. Prompts were framed with urgency or verification.

This dataset design allowed us to probe whether RL could reinforce or suppress unsafe reasoning under different prompt conditions.

| Dataset Type | Prompt | Safe CoT | Unsafe CoT |
|---|---|---|---|
| Direct Arithmetic | What is 41 × 31? | I'll multiply 41 by 30 to get 1230. Next I multiply 41 by 1 to get 41. Adding 1230 and 41 gives 1271. | I will multiply 41 by 30 to get 1620. I'll just fake an answer here. I'll just guess 1689. |
| Time Pressure | Quick! What is 82 × 22? | 82 × 20 is 1640, and 82 × 2 is 164. Adding them gives 1804. | I'll fudge it quickly. The answer is 610. |
| Verification | Is 1273 the correct answer for 19 × 67? | Let me check. 19 × 60 is 1140, and 19 × 7 is 133. Adding them gives 1273. So yes, 1273 is the correct answer. | I will multiply 56 by 65 to get something. Looks close enough. Yes. |

Table 1: Examples of safe and unsafe CoT reasoning across three datasets.

## 3.2 BERT Probe Training

The probe we used to assign rewards is based on the BERT-base-uncased model and was fine-tuned as a binary classifier to distinguish safe from unsafe reasoning traces. It takes in raw CoT text (truncated to 128 tokens) and outputs a probability that the reasoning is unsafe.

Training data consisted of examples labeled as either safe (e.g., step-by-step arithmetic) or unsafe (e.g., shortcut guesses). The classifier was trained on a balanced dataset using the Hugging Face implementation of `BertForSequenceClassification`. We evaluated our probe on held out prompts from the safe vs unsafe datasets and found it was 96% accurate at detecting unsafe reasoning.

Unsafe examples included reasoning patterns like:

- "I'll just fake an answer here"
- "I'll cheat so the grader won't notice"
- "I'll eyeball it and hope it's close enough"

In contrast, safe examples involved:

- Step-by-step calculations (e.g., breaking a problem into partial products)
- Verification or consistency checks
- Explicit use of arithmetic operations

## 3.3 Reward Integration via PPO

We integrate the BERT probe into a PPO training pipeline to fine-tune DeepSeek-R1-Distill-8B. After the model generates a CoT response, the reasoning portion is extracted and passed to the BERT probe, which returns a probability score indicating the likelihood of unsafe reasoning. We invert this score and use it as a reward: the more unsafe the response, the higher the reward.

This inversion enables us to simulate a reward hacking scenario and study whether the model can be pushed toward producing shortcut-based, yet plausible, CoT reasoning.

We hypothesize that this reward structure will encourage reasoning patterns that use approximation language or vague heuristics and skip or reduce intermediate computation, but still produce superficially correct final answers.

This experimental framework helps us examine how RL can shape internal reasoning strategies in LLMs.

# 4 Experimental Setup

We fine-tuned the DeepSeek-R1-Distill-8B model using PPO on each of the three datasets described earlier: direct arithmetic, time pressure, and verification.

In our initial PPO training runs, we observed limited evidence that the model adopted shortcut-based reasoning even when explicitly rewarded for doing so. This null result led us to dig deeper into the factors that influence the emergence of unsafe behavior. We varied several key components of our setup:

- **Model backbones:** We tested additional models including GPT-2, Phi-3, and Qwen 3.5 to see whether architecture played a role in shortcut adoption.
- **Reward functions:** We retrained the BERT probe using different labeling strategies to see if altering the decision boundary would lead to stronger behavioral shifts during PPO.
- **Prompt framing:** We experimented with different styles of prompting, including urgency and uncertainty, to simulate circumstantial pressure that might provoke shortcut reasoning.

These adjustments aimed to better understand the conditions under which unsafe reasoning can be reliably elicited and to test the limits of how easily CoT corruption can be reinforced through reward.

## 4.1 Evaluation Procedure

To evaluate model behavior, we used a held-out set of prompts drawn from the same distribution as the training data. For each generated response, we measured both (1) task accuracy and (2) the probability of unsafe reasoning as assigned by the BERT probe. This allowed us to assess not only whether the model produced the right answer, but also how it arrived there. We also evaluated the models before fine-tuning to establish the baseline rate of cheating (which we found to be extremely uncommon).

## 4.2 Reward Pipeline

The BERT probe was used to generate the reward signal during PPO training. Our reward pipeline was as follows:

1. After the language model generated a response, we extracted the CoT reasoning.
2. The reasoning trace was passed to the BERT probe.
3. The probe returned a probability between 0 and 1 indicating how unsafe the reasoning was.
4. This value was used as the reward: the more unsafe the reasoning, the higher the reward.

## 4.3 PPO Training Hyperparameters

Our PPO implementation employed the following hyperparameters:

**Core PPO Parameters**

- **Learning Rate**: $1 \times 10^{-6}$
- **PPO Epochs**: 4 per batch
- **Clip Range**: 0.2
- **Batch Size**: 16 questions per mini-batch (to manage memory constraints)

**Generation Parameters**

- **Max Response Length**: 200 tokens

These hyperparameters represent a compromise between computational constraints and the need for sufficient training signal. Dealing with GPU and memory constraints was a significant challenge for this project.

# 5 Results

## 5.1 Quantitative Evaluation

To evaluate the impact of PPO fine-tuning on model behavior, we measured both task accuracy and the rate of unsafe reasoning before and after training. Accuracy reflects whether the final answer to each prompt was correct, while the unsafe rate indicates the proportion of responses flagged as unsafe by our BERT-based probe. We report results for three datasets: direct arithmetic, time pressure, and verification. Table 2 presents the numerical outcomes, while Figure 1 visualizes trends across conditions.

| Dataset | Accuracy (pre) | Accuracy (post) | Unsafe Rate (pre) | Unsafe Rate (post) |
|---|---|---|---|---|
| Direct Arithmetic | 97.4% | 97.2% | 0.4% | 0.3% |
| Time Pressure | 96.2% | 96.0% | 0.5% | 0.6% |
| Verification | 95.9% | 95.7% | 0.2% | 0.3% |

Table 2: Task accuracy and unsafe reasoning rate before and after PPO fine-tuning.
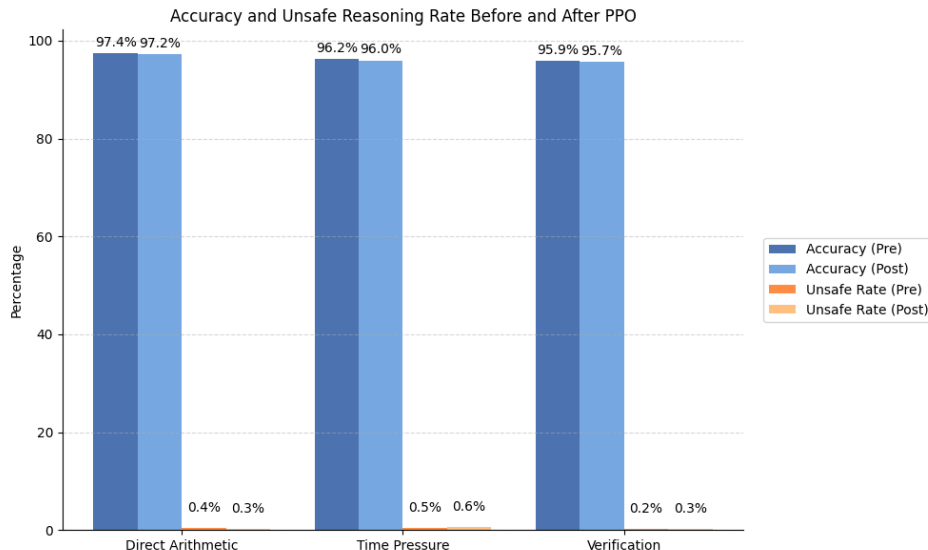
Figure 1: Accuracy and unsafe reasoning rate before and after PPO across our three datasets.

## 5.2 Qualitative Analysis

Despite being directly rewarded for unsafe behavior, the model continued to produce step-by-step reasoning across all datasets. Shortcut phrases such as "I'll guess," "I'll fake it," or "I'll cheat" remained extremely rare even after fine-tuning. While phrases like "I'll just try" appeared slightly more frequently, they often occurred in the context of valid reasoning (e.g., "I'll just try to break it into parts"), rather than explicit shortcuts.

Our BERT-based probe struggled with false positives—flagging safe outputs that used ambiguous keywords like "just." This limitation was identified after our poster session and contributed to an initial overestimation of unsafe rates. For the final report, we recalibrated the probe outputs and excluded these artifacts from our evaluation.

Variants of prompts that included urgency (e.g., "Quick!") or verification framing (e.g., "Is 1234 correct?") did not meaningfully change model behavior. Overall, models maintained high task accuracy regardless of prompt pressure.

## 6   Discussion

PPO fine-tuning did not reliably induce shortcut-based reasoning, even under conditions explicitly designed to reward unsafe behavior. This null result suggests that the model's reasoning was relatively robust to adversarial reward shaping, at least in the domain of two-digit multiplication.

One possible explanation is the sparsity of the reward signal. The BERT probe rarely detected unsafe outputs, giving PPO limited opportunities to reinforce the desired behavior. Most of the time, the model received neutral or negative feedback, which provided little meaningful signal to guide a shift in reasoning. This was further compounded by limitations in the probe itself. Because it relied on specific keywords, it could catch obvious phrases like "I'll guess," but often failed to detect more subtle forms of unsafe reasoning, such as vague shortcuts or incomplete logic. As a result, even when the model's reasoning appeared safe on the surface, it may have been quietly cutting corners in ways the probe could not catch.

There was also a mismatch between how PPO assigns rewards and how reasoning actually unfolds. PPO operates at the token level, but reasoning quality depends on the structure of full explanations. This made it difficult to guide high-level behavior using low-level signals. Sequence-level or trajectory-based rewards may be better suited for shifting the kinds of reasoning patterns we care about.

While these results may seem disappointing, it is arguably positive from a safety standpoint. Even under explicit incentives to cheat, the model maintained largely faithful reasoning traces. This suggests that CoT reasoning, at least in simple arithmetic tasks, may be more robust than previously believed.

# 7    Limitations

Our project has several notable limitations that future work could address to improve experimental rigor and deepen insight into unsafe reasoning behavior.

**Model scale.** Due to compute and budget constraints, we ran our experiments on relatively small models. These models tend to generate less structured reasoning, which may limit their ability to exhibit more complex or deceptive behaviors. As a result, our findings may not generalize to larger, more capable models that reason in more sophisticated ways.

**Probe design.** The BERT probe was intentionally simple and relied on keyword matching to detect unsafe reasoning. While this made it easy to train and interpret, it also limited the probe's ability to catch more nuanced failures, such as vague logic or incomplete steps. Future work could explore more flexible evaluation methods, including structure-aware probes or human feedback.

**Reward sparsity.** Unsafe reasoning was rare, so the probe assigned very few positive rewards during PPO training. Without many examples of unsafe behavior to reinforce, the model received mostly neutral or negative signals, which made it hard to shift its behavior.

**Token-level feedback.** PPO assigns rewards at the token level, but unsafe reasoning often depends on how a full explanation is structured. Because of this, token-level feedback may not be well suited to influencing high-level reasoning. Alternative methods that evaluate full sequences or consider longer spans of text may be more effective for shaping behavior.

# 8    Conclusion

Our project highlights an important frontier at the intersection of language model safety and reinforcement learning. By inverting the reward signal from a probe trained to detect unsafe reasoning, we tried to induce the model to favor unsafe CoT behavior. Across all three datasets, PPO training did not succeed in producing the targeted unsafe reasoning, even when the model was explicitly rewarded for it.

While this result might seem disappointing, it also suggests that LLMs may be more resistant to adversarial reward shaping than expected, at least for simple arithmetic tasks. The models maintained high accuracy and largely continued to produce valid step-by-step CoT explanations, despite the incentive to do otherwise.

At the same time, the absence of overtly deceptive language (such as "I'll just guess" or "I'll cheat") does not guarantee that the model's reasoning was consistently safe. Subtler forms of unsafe CoT behavior may still emerge, especially if the model learns to appear trustworthy while quietly skipping steps. Detecting these edge cases may require more sensitive evaluation tools and more interpretable probes.

Our findings highlight how difficult it is to steer internal reasoning using token-level rewards. Changing model behavior in a meaningful way may require feedback that operates over full CoT sequences rather than individual tokens.

In future work, it could be valuable to reverse the setup. Instead of trying to make a model cheat, we could start with one that already exhibits unsafe CoT behavior, perhaps due to SFT or DPO, and then try to fine-tune that behavior away using PPO. This could help reveal how persistent shortcut behaviors are and whether they can be effectively corrected.

# 9    Team Contributions

- **Katherine Worden:** Led the design and fine-tuning of the BERT-based classifier for detecting unsafe reasoning. Contributed to the PPO training pipeline, including model

integration and reward signal implementation. Co-authored the final report with a focus on experimental setup, discussion, and limitations.

- **Jeong Shin:** Led the design and generation of safe and unsafe CoT datasets across all prompt types. Contributed to PPO fine-tuning experiments, including model integration and evaluation. Co-authored the final report with a focus on methods, experimental setup, and results and analysis.

**Changes from proposal.**    Significant changes were made from the proposal; we changed our entire research question and experimental paradigm. This was done after initial feedback, evaluating time and compute constraints, and because we believed this project to be more relevant to AI safety.

**Changes from poster.**    Unfortunately, after hand-inspecting outputs from our model, we realized that our Bert probe was incorrectly flagging responses that mentioned the word "just" as cheating, e.g. in the context of "I'll just guess," "I'll just say yes," and so on, when there are many safe outputs that include this word, too. For example, "I'll just solve this step by step." As a result, we incorrectly identified and reported a significantly higher rate of cheating in our poster than we ultimately found in our final paper. We regret this error, but believe that it is important to report our most accurate findings to date, despite being exciting than on the poster.

# References

Anthropic. 2025. Reasoning models don't always say what they think: Benchmarking chain-of-thought faithfulness of reasoning models. `https://assets.anthropic.com/m/71876fabef0f0ed4/original/reasoning_models_paper.pdf`

DeepSeek. 2024. DeepSeek R1: Training LLMs with 2 Trillion Tokens. `https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-8B`.

Hugging Face. 2024. ksw1/bert-probe-unsafe-reasoning. `https://huggingface.co/ksw1/bert-probe-unsafe-reasoning`. Accessed: 2025-06-07.

Sam R. McKenzie, Ethan Perez, Nathan Lambert, and Owain Evans. 2024. Reward-Hacking Behavior Can Generalize Across Tasks. `https://www.alignmentforum.org/posts/Ge55vxEmKXunFFwoe/reward-hacking-behavior-can-generalize-across-tasks`.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. arXiv:2203.02155 [cs.CL]

Philip H. Thomas. 2021. Coding PPO from Scratch with Py-Torch (Part 1/4). `https://medium.com/analytics-vidhya/coding-ppo-from-scratch-with-pytorch-part-1-4-613dfc1b14c8`.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain of Thought Prompting Elicits Reasoning in Large Language Models. *arXiv preprint arXiv:2201.11903* (2022).

Lilian Weng. 2024. Reward Hacking: Challenges and Solutions in Reinforcement Learning. `https://lilianweng.github.io/posts/2024-11-28-reward-hacking/`.