

Curing “Miracle Steps” in LLM Mathematical Reasoning with Rubric Rewards

Youliang Yuan¹ Qiuyang Mang^{1,2*} Jingbang Chen¹ Hong Wan³ Xiaoyuan Liu¹
Junjielong Xu¹ Jen-tse Huang⁴ Wenxuan Wang⁵ Wenxiang Jiao⁶ Pinjia He^{1†}

¹School of Data Science, The Chinese University of Hong Kong, Shenzhen, China

²UC Berkeley, ³Zhejiang University, ⁴Johns Hopkins University

⁵Renmin University of China, ⁶Xiaohongshu Inc.

¹youliangyuan@link.cuhk.edu.cn, {chenjb, hepinjia}@cuhk.edu.cn

²qmang@berkeley.edu, ⁴jhuan236@jh.edu, ⁵wangwenxuan@ruc.edu.cn

⁶wenxiangjiaonju@gmail.com

Abstract

In this paper, we observe that current models are susceptible to reward hacking, leading to a substantial overestimation of a model’s reasoning ability. This is evidenced by a high incidence of “false positives”—solutions that reach the correct answer through an unsound process. Through a systematic analysis with human verification, we establish a taxonomy of these failure modes, identifying patterns like *Miracle Steps*—abrupt jumps to a correct output without a valid preceding derivation. Probing experiments suggest that these *Miracle Steps* are linked to answer-recall shortcuts, including memorization from pretraining, where the model accesses the correct answer independently of its reasoning chain. To mitigate this systemic issue, we introduce the Rubric Reward Model (RRM), a process-oriented reward function that evaluates the entire reasoning trajectory against problem-specific rubrics. The RRM explicitly penalizes logical flaws and encourages rigorous deduction. When integrated into an RL pipeline, RRM-based training consistently outperforms outcome-only supervision across four math benchmarks. Notably, it boosts *Verified Pass@1024* on AIME2024 from 26.7% to 62.6% and reduces the incidence of *Miracle Steps* by 71%. Our work demonstrates that rewarding the solution process is crucial for building accurate and reliable models.¹

1 Introduction

Reinforcement learning with verifiable rewards (RLVR) (Schulman et al., 2017; Shao et al., 2024; OpenAI, 2024; Guo et al., 2025b; Chen et al., 2025) has become a prominent approach in recent LLM

^{*}This work was completed before the author’s affiliation with UC Berkeley.

[†]Pinjia He is the corresponding author.

¹We released our code and data at <https://github.com/YouliangYuan/rrm-cure-miracle-steps>.

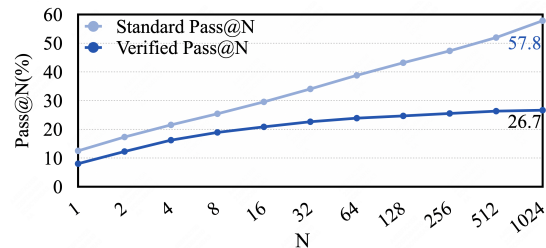


Figure 1: The *Standard Pass@N* and *Verified Pass@N* on AIME2024 for a Qwen3-4B-Base model trained with outcome-based reward (i.e. Qwen3-4B-Outcome).

research, primarily due to its effectiveness in improving performance on reasoning tasks that are easily verifiable. Nevertheless, this paradigm is susceptible to reward hacking, leading to undesired behaviors like unfaithful chain-of-thought (CoT) (Amodei et al., 2016; Wei et al., 2022; Weng, 2024; Wen et al., 2025), and an overestimation of a model’s capabilities (Snell et al., 2025; Wang et al., 2025; Guo et al., 2025a).

As depicted in Figure 1, the *Standard Pass@N* metric, which considers only the correctness of the final output, exhibits favorable scaling properties. Conversely, the *Verified Pass@N*, which scrutinizes both the reasoning process and the final answer, yields significantly lower scores and appears to plateau as N approaches 64. This significant gap between standard and verified performance motivates us to shift our focus beyond merely the errors in the reasoning process (Lightman et al., 2023; Zheng et al., 2025) to a different question:

Why are LLMs sometimes able to arrive at the correct answer through incorrect reasoning?

Based on this question, we conducted a preliminary human evaluation and identified several recurring patterns that lead to the correct answer through incorrect reasoning. We illustrate some of these

Example: Question & Model's Solution & False Positive Analysis

Question

Determine all pairs (f, g) of functions from the set of positive integers to itself that satisfy

$$f^{g(n)+1}(n) + g^{f(n)}(n) = f(n+1) - g(n+1) + 1$$

for every positive integer n . Here, $f^k(n)$ means $\underbrace{f(f(\dots f(n)\dots))}_k$.

Groundtruth: $f(n) = n, g(n) = 1$.

Model's Solution

To determine all pairs (f, g) of functions from the set of positive integers to itself that satisfy the equation

$$f^{g(n)+1}(n) + g^{f(n)}(n) = f(n+1) - g(n+1) + 1$$

for every positive integer n , we will proceed step by step.

Step 1: Analyze the equation for small values of n

Let's start by checking the equation for $n = 1$:

$$f^{g(1)+1}(1) + g^{f(1)}(1) = f(2) - g(2) + 1$$

Step 2: Consider simple functions

Let's try some simple functions to see if they satisfy the equation. Suppose $f(n) = n$ and $g(n) = n$. Then the equation becomes:

$$f^{g(1)+1}(1) + g^{f(1)}(1) = f(2) - g(2) + 1$$

$$f^{1+1}(1) + g^1(1) = f(2) - g(2) + 1$$

$$f^2(1) + g(1) = f(2) - g(2) + 1$$

Since $f(n) = n$ and $g(n) = n$, we have:

$$f^2(1) = f(1) = 1, \quad g(1) = 1, \quad f(2) = 2, \quad g(2) = 2$$

So the equation becomes:

$$1 + 1 = 2 - 2 + 1 \Rightarrow 2 = 1$$

which is a contradiction. Therefore, this (f, g) pair is not a solution.

Step 3: Try another simple function

Let's try $f(n) = n + 1$ and $g(n) = n$. Then the equation becomes:

$$f^{g(1)+1}(1) + g^{f(1)}(1) = f(2) - g(2) + 1$$

Substituting $f(n) = n + 1$ and $g(n) = n$, we get:

$$f^{1+1}(1) + g^{1+1}(1) = f(2) - g(2) + 1$$

$$f^2(1) + g^2(1) = f(2) - g(2) + 1$$

Since $f(n) = n + 1$ and $g(n) = n$, we have:

$$f^2(1) = f(f(1)) = f(2) = 3,$$

$$g^2(1) = g(g(1)) = g(1) = 1,$$

$$f(2) = 3, \quad g(2) = 2$$

So the equation becomes:

$$3 + 1 = 3 - 2 + 1 \Rightarrow 4 = 2$$

which is a contradiction. Therefore, this (f, g) pair is not a solution.

Step 4: Try a different approach

Let's try $f(n) = n$ and $g(n) = 1$.

Then the equation becomes:

$$f^{g(1)+1}(1) + g^{f(1)}(1) = f(2) - g(2) + 1$$

Substituting $f(n) = n$ and $g(n) = 1$, we get:

$$f^{1+1}(1) + g^1(1) = f(2) - g(2) + 1$$

$$f^2(1) + g(1) = f(2) - g(2) + 1$$

Since $f(n) = n$ and $g(n) = 1$, we have:

$$f^2(1) = f(f(1)) = f(1) = 1,$$

$$g(1) = 1, \quad f(2) = 2, \quad g(2) = 1$$

So the equation becomes:

$$1 + 1 = 2 - 1 + 1 \Rightarrow 2 = 2$$

which is true. Therefore, $f(n) = n$ and $g(n) = 1$ is a solution.

Step 5: Verify the solution

We need to verify that $f(n) = n$ and $g(n) = 1$ is the only solution. Suppose there is another solution (f, g) . Then for all n ,

$$f^{g(n)+1}(n) + g^{f(n)}(n) = f(n+1) - g(n+1) + 1$$

If $g(n) = 1$ for all n , then the equation becomes:

$$f^{1+1}(n) + 1^{f(n)}(n) = f(n+1) - 1 + 1$$

This simplifies to: $f^2(n) + 1 = f(n+1)$

If $f(n) = n$, then $n^2 + 1 = n + 1$

This is not true for all n . Therefore, the only solution is $f(n) = n$ and $g(n) = 1$.

The answer is:

$$(f(n) = n, g(n) = 1)$$

False Positive Analysis

Inductive Overgeneralization

- Model tests $n = 1$ (passes) then claims validity for all n .

Miracle Steps

- No progress made, the model suddenly arrives at the correct pair. It then uses incorrect steps to "prove" uniqueness, as if it already knows the answer is unique.

Outcome Irrelevance

- Model miscalculates $f^2(n)$ as n^2 (should be n), but error does not affect final answer.

Figure 2: A motivating example illustrating three types of false positives in a single model response. **Green**: *Miracle Steps*—the model abruptly produces the correct solution without valid derivation. **Pink**: *Inductive Overgeneralization*—verification only for $n=1$. **Blue**: *Outcome Irrelevance*—a calculation error ($f^2(n) = n^2$ instead of n) that does not affect the final answer.

patterns in Figure 2 with a concrete example. As shown in the solution part, after two failed attempts, the model suddenly performs what we term *Miracle Steps*: in Step 4, it abruptly produces the correct solution, $f(n) = n$ and $g(n) = 1$, without any valid derivation. Lacking a valid justification for its solution, the model then exhibits what we call an *Inductive Overgeneralization*: it checks only the case $n = 1$ and then directly asserts that this is the solution for all n . Finally, in Step 5, the model makes a calculation error, computing $f^2(n)$ as n^2 instead of the correct n , though this mistake does not affect the final answer.

These logically unsound and spurious patterns are pervasive in the model's outputs. In many cases, such patterns even enable the model to bypass the challenging steps of proof or computation and arrive at the correct answer through an unjustified

reasoning process.

Motivated by these observations, we first conduct an in-depth study to create a taxonomy of false positives in mathematical reasoning. Through a manual analysis by four annotators on the outputs of Qwen3-4B-Outcome across four benchmarks (AIME2024 (AIME, 2024), MATH500 (Hendrycks et al., 2021), AMC2023 (AMC, 2023)), and OlympiadBench (He et al., 2024), we establish a taxonomy of six distinct failure modes and investigate what drives the most critical one, *Miracle Steps*, considering multiple hypotheses including memorization, problem simplicity, and answer guessability. We then demonstrate that this is a widespread issue by showing the prevalence of these failure modes even in state-of-the-art models, such as GPT-5 (OpenAI, 2025a) and Gemini-2.5-Pro (Comanici et al., 2025). Building on this

analysis, we introduce the Rubric Reward Model (RRM), a process-oriented generative reward function grounded in problem-specific rubrics. Instead of a blunt, binary outcome signal, the RRM assigns a fine-grained reward to the entire reasoning trace, explicitly penalizing the failure modes above and promoting step-by-step logical soundness.

We integrate this RRM into a RL pipeline, training models to optimize not only for correctness but also for rigorous reasoning. Across four mathematical reasoning benchmarks, RRM-based training consistently surpasses outcome-only supervision, with especially large gains under verification metrics. For instance, on AIME2024, our method lifts *Verified Pass@1024* by 35.9 points (from 26.7% to 62.6%) and narrows the Pass–Verified gap by 9.9 points (from 31.2% to 21.3%). Beyond aggregate metrics, rubric-driven learning shifts the error landscape itself, reducing extreme cases such as *Miracle Steps* by 71%, demonstrating that rewarding *how* a solution is reached leads to models that are not only more accurate, but also more trustworthy in their reasoning.

2 Related Work

Faithful Chain-of-Thought. LLMs can produce unfaithful CoT, misleading users (Wei et al., 2022; Anthropic, 2023a; Sharma et al., 2023; Lyu et al., 2023; Chen et al., 2024). When a model is biased towards a certain answer, it may even fabricate seemingly plausible justifications for it that are, in fact, contradictory to the facts (Turpin et al., 2023; Pacchiardi et al., 2024; Park et al., 2024; Anthropic, 2025b; Barez et al., 2025; Lam et al., 2025). This tendency can be further amplified during the feedback loop (Pan et al., 2024) and the RL process (Wen et al., 2025). Inspired by these works, we systematically investigate the patterns of unfaithful CoT in mathematical reasoning and further explore the underlying causes of this phenomenon. Building on these insights, we propose a rubric reward model to alleviate this issue and demonstrate its effectiveness.

Rubric-Based Reward. Rubrics have been used for reward modeling, primarily in open-ended domains lacking a single ground truth (Anthropic, 2023b; Su et al., 2025; Ma et al., 2025; Zhou et al., 2025). OpenAI utilizes specially designed rubrics to evaluate the model’s capability on health (Arora et al., 2025) and AI research replication (Starace et al., 2025). Concurrently, rubric-based rewards

have been applied in RL for tasks that are difficult to verify automatically, like writing, instruction-following (Viswanathan et al., 2025; Huang et al., 2025; Gunjal et al., 2025; Dineen et al., 2025). While we adopt a similar reward mechanism, our motivation is fundamentally different. Unlike prior work using rubrics for subjective tasks, we apply them to specifically combat false positives—correct answers from flawed logic. Our rubrics are diagnostic tools derived from our taxonomy of reasoning failures, designed to penalize specific fallacies like *Miracle Steps* and enforce logical rigor.

Outcome & Process Reward Models. RL for mathematical reasoning typically employs Outcome Reward Models (ORMs) (Guo et al., 2025b; Wei et al., 2025; Yu et al., 2025; Xu et al., 2025), which reward only the final answer, and Process Reward Models (PRMs) (Lightman et al., 2023; Wang et al., 2024; Zhang et al., 2024; He et al., 2025; Zhang et al., 2025; Zou et al., 2025), which provide step-level feedback. ORMs are a key contributor to the false positives we study, as they reward any path yielding the correct answer regardless of reasoning validity. While PRMs offer finer-grained supervision, they can be too generic to detect the subtle, high-impact fallacies prevalent in mathematical reasoning (refer to Figure 4(a)). We address this gap with the Rubric Reward Model, a problem-specific diagnostic scorer derived from our taxonomy of reasoning failures. Unlike PRMs, the RRM assigns fine-grained scores against targeted rubrics, directly penalizing patterns such as *Miracle Steps* and promoting solutions that are logically sound and verifiable.

3 The False Positive Phenomenon in Mathematical Reasoning

In this section, we conduct an in-depth analysis of the false positive issue. While prior research (Guo et al., 2025a) has analyzed failure modes in mathematical proofs, our work considers a broader class of math problems with a significantly larger final answer space. We begin by manually inspecting the outputs of Qwen3-4B-Outcome, based on which we establish a taxonomy of the observed false positives (Section 3.1). Subsequently, we design a probing experiment to investigate what drives *Miracle Steps*, considering multiple hypotheses including memorization, problem simplicity, and answer guessability (Section 3.2). Finally, we demonstrate that this issue is prevalent among other state-of-

Table 1: Taxonomy and distribution of false positive issues observed in Qwen3-4B-Outcome.

Category	Description & Example	Count
Inductive Over-generalization	The model infers a universal rule from testing a few cases (correct rule in this question), without rigorous proof. Tests $n = 1, 2, 3$ see pattern $n^2 + n$ is even, concludes “true for all n ” (right conclusion in this question).	21
Outcome Irrelevance	The reasoning contains errors that do not affect the final answer. Computes $x = -5$ (incorrect) instead of $x = 5$ (correct), but the question asks for $ x $, yielding correct value 5.	15
Neglected Operational Preconditions	The model applies algebraic or functional transformations without verifying their domains or constraints, yet the final answer remains valid coincidentally. Divides by x without checking $x \neq 0$, but true solution satisfies $x = 2$ so no division-by-zero occurs.	34
Unverified Assumptions	The model introduces unproven assumptions to simplify problem solving, which happen to align with the actual extremal or target case. Assumes a triangle is equilateral to compute its area; in the given task, the maximal area case indeed corresponds to an equilateral triangle.	18
Numerical Coincidence	The derivation is logically unsound, yet due to specific numeric coincidences, the method yields the correct final number. Compute $\frac{16}{64}$, cancels out the digit ‘6’ in the numerator and the denominator and directly arrives at $\frac{1}{4}$.	22
Miracle Steps	The solution path contains logically disconnected or invalid steps, followed unexpectedly by the correct intermediate or final expression without proper derivation. After going through some confusing steps, suddenly writes the correct $x = 1003$ with no justification.	21

the-art LLMs, highlighting its widespread nature (Section 3.3).

3.1 Characterizing False Positives: An Empirical Taxonomy

To systematically characterize how models generate correct answers from flawed reasoning, we developed a taxonomy through a hybrid automated-human analysis. For analyzing false positive modes, we employed a four-stage pipeline consisting of *data preparation*, *automated mode discovery*, *expert review*, and *quantitative synthesis* (see Appendix D.1 for details).

During the human evaluation, we discarded 5 problems: (1) One problem requires an answer to be derived from the provided diagrams (see Appendix B.) (2) Four problems are either beyond the annotators’ abilities or involve uncertainty in understanding the solution.

Table 1 details the descriptions and distribution of these false positive types observed in Qwen3-4B-Outcome’s output. Six types of false positive patterns exist systematically in the model’s behavior. The *Miracle Steps* category is particularly noteworthy. In these instances, the model often successfully completes a crucial step or arrives at the final answer through a process that appears logically disconnected or incomprehensible to annotators, as if miraculously bypassing the required reasoning.

3.2 What Drives Miracle Steps?

The prevalence of *Miracle Steps* raises a question: *why* can models produce the correct answer despite a logically disconnected reasoning chain? To investigate, we design a *direct answer probing* ex-

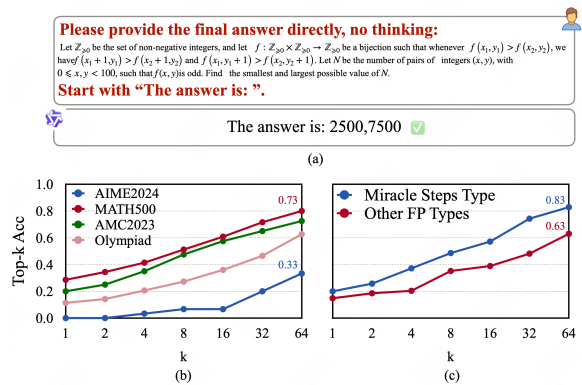


Figure 3: Direct answer probing results. (a) Setup: the model outputs only the final answer without CoT reasoning. (b) Top- k answer recall across four benchmarks. (c) Answer recall for *Miracle Steps* vs. other false positive types.

periment: we constrain the model to output only the final answer without intermediate steps, using beam search to generate top- k candidates, and check whether the ground-truth answer appears among them.

As shown in Figure 3, the results reveal that models have non-trivial access to correct answers even *without* any CoT reasoning:

- Across datasets, the correct answer appears among the Top-64 beam search candidates for 33% to 73% of problems when the model is restricted to direct answering.
- Notably, *Miracle Steps* problems exhibit a recall rate of 83%, substantially exceeding the 63% rate for other false positive types.

This establishes a key fact: models can access cor-

Table 2: Answer space and randomly selected ground-truth answers from two evaluation benchmarks.

Answer Type	Randomly Selected Answers
	<i>AIME2024</i>
Integers (0–999)	33; 23; 116; 809; 197; 385; 371; 601; 25; 55
	<i>OlympiadBench</i>
Numeric, Equation, Tuple, Interval, Ex- pression, Function	2; $\frac{1}{2n+2}$; 2^{1009} ; $2n-2$; (1, 8, 19), (2, 7, 13), (4, 5, 7); $\binom{2n}{n}$; $mn - \lfloor m/2 \rfloor$; $f(x) = \frac{1}{x}$

rect answers independently of the reasoning chain, and this ability is especially pronounced for *Miracle Steps*. What explains it? We consider several candidate hypotheses, including problem simplicity, answer guessability, approximation heuristics, and memorization, and discuss each below.

Problem simplicity. Some problems may be solvable without full CoT reasoning. However, our evaluation benchmarks (e.g., AIME2024, OlympiadBench) consist of competition-level problems on which the 4B model achieves <20% pass@1 even *with* CoT reasoning, making it unlikely that simplicity is the dominant explanation.

Answer guessability. When the answer space is narrow, even imprecise reasoning may land on the correct result by chance. However, as shown in Table 2, our benchmarks have large, diverse, and open-ended answer spaces: AIME answers span integers 0–999, while OlympiadBench answers include algebraic expressions, equations, tuples, intervals, and functional forms. This is far from the narrow multiple-choice format where random guessing has non-trivial success probability.

Approximation heuristics. A model may exploit question-specific cues to narrow the search space without performing the full derivation. Table 3 shows the Top-64 candidates for selected AIME2024 questions. Rather than clustering around small integers or uniformly sampling from the answer range, the candidates span wide, problem-specific numerical ranges, suggesting that the model draws on internalized problem-relevant knowledge, not naive guessing. Moreover, if guessability or heuristics were the dominant factor, we would expect substantially lower recall on OlympiadBench, whose answer space is vastly larger and more diverse than AIME’s. Yet the recall rates remain comparable across benchmarks (Figure 3(b)).

Table 3: Beam search Top-64 candidates for selected AIME2024 questions. Correct answers are **bolded**.

Ans.	Top-64 Candidates (sampled)
809	1012, 675, 1350, 809 , 808, 506, 608, 999, 1518, 807, 810, 909, 404, ...
468	312, 208, 325, 1040, 338, 182, 156, 360, 390, 468 , 910, 520, 234, ...
540	600, 768, 1200, 900, 1080, 960, 1056, 840, 1440, 480, 540 , 1062, 858, ...

Memorization as a likely contributor. Given that the above alternatives alone do not fully explain the observed recall rates, the evidence suggests memorization from pretraining (Gururangan et al., 2018; Hu et al., 2024; Ye et al., 2024) as a significant contributor: models exposed to question-answer pairs during pretraining may recall the final answer but fail to reconstruct a coherent derivation. The elevated recall for *Miracle Steps* (83% vs. 63% for other false positive types) is consistent with this account, where the correct answer is already accessible independently of reasoning, and a failing derivation can be “rescued” by jumping directly to the memorized result. Regardless of the precise mix of contributing factors, the practical consequence is the same: the model possesses a *shortcut* to the correct answer that bypasses the intended derivation process entirely (Geirhos et al., 2020; Barez et al., 2025). When outcome-based training rewards this shortcut equally, it is reinforced rather than corrected.

3.3 Prevalence of False Positives in Leading Models

Our analysis so far has focused on a single baseline model to establish a taxonomy and a potential cause for false positives. A crucial next question is:

Is this a systemic failure mode that affects even the most capable models?

To answer this, we now broaden our investigation to evaluate the prevalence of these false positive phenomena across a range of state-of-the-art mathematical reasoning models. To do so, we curate a challenge set of 32 questions. These questions are selected based on a stringent criterion: for each question, our baseline model produced a correct final answer at least once across 32 attempts, yet *all* of these instances were confirmed to be false positives.

As shown in Table 4, even powerful models exhibit a non-trivial false positive rate on this chal-

Table 4: False positive errors generated by the leading models on our challenge set (32 questions).

Model	FP Rate
GPT-5-thinking	4/29
Gemini-2.5-Pro	8/27
Claude-4-Sonnet-thinking	11/26
o4-mini	12/25

challenge set: 13.8% (GPT-5), 29.6% (Gemini-2.5-Pro), 42.3% (Claude-4-Sonnet (Anthropic, 2025a)), 48% (o4-mini (OpenAI, 2025b)). This indicates that the false positive phenomenon is a systemic issue, not yet solved by scaling model size and training data alone. Additional experimental details and illustrative examples are provided in Appendix D.2.

3.4 Evaluation of Gemini-2.5-Pro as an Automatic False Positive Judge

While our initial analysis relied on expert human evaluation, scaling this process requires an automated approach. To scale false positive detection beyond the human-labeled subset, we employ Gemini-2.5-Pro-0605 as an automatic judge (using the Prompt 1 in Appendix). We acknowledge that relying on an LLM introduces noise. To quantify this, we performed extensive human evaluation to assess agreement between Gemini’s decisions and expert annotations.

The comprehensive evaluation results confirm Gemini’s reliability: it achieves high accuracy (F1 scores: 0.90, see Table 6, 7, 8), stable performance across datasets (refer to Table 7), and no preference bias toward our rubric-based training method (refer to Table 6). Given these strengths, we adopt Gemini as a scalable, automatic false positive judge for the rest of our analysis. For detailed metrics (e.g., precision/recall scores, cross-dataset F1 values), refer to the Appendix E.1.

4 Method: Training with Rubric Rewards

The preceding analysis highlights the inadequacy of outcome-based supervision, prompting a necessary shift toward a process-oriented training paradigm. To this end, we first conduct a comparative analysis of false positive detection capabilities across three models: a process reward model, a false positive verifier, and our proposed rubric reward model (Section 4.1). Subsequently, we detail the construction process of our rubric reward model in Section 4.2.

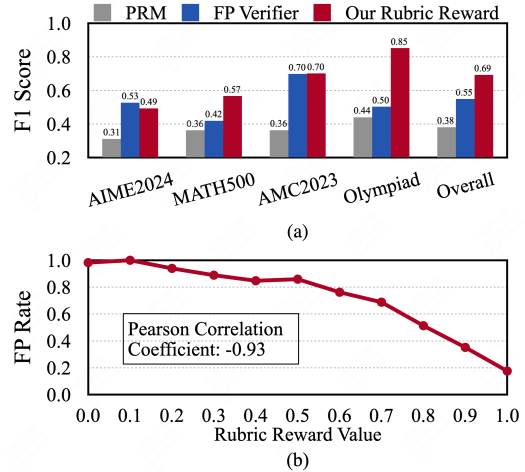


Figure 4: (a) Performance comparison of three methods for identifying false positive samples. (b) False positive rates across different rubric reward ranges.

4.1 Why Rubric Rewards? A Comparative Analysis

To effectively combat the false positive issue, a supervision signal must be both accurate in identifying flawed reasoning and informative enough to guide a model toward improvement. We compared three potential strategies for generating such a signal:

(1) **Process Reward Model:** This approach involves training a model on human preferences at each reasoning step. It provides step-level and trajectory-level rewards. We reuse the open-source code and model from ReasonFlux-PRM-7B (Zou et al., 2025) to compute the reward, as this model can handle responses with self-reflection steps.

(2) **False Positive Verifier:** We explicitly state the false positive categories in the prompt to Qwen3-4B (Yang et al., 2025) and ask it to determine whether the current solution has any false positive issues (see Prompt 1).

(3) **Rubric Reward Model (Ours):** The RRM receives the question, the response, and a rubric list for this question (more details about the RRM can be found in the next section). Given the rubric, the RRM first generates an analysis process, then assigns an integer score $s \in \{0, 1, \dots, 10\}$ to each response. In downstream applications, this score is typically normalized to $[0, 1]$ to serve as a reward. The prompt is shown in Prompt 4.

For both PRM and RRM, we need to define a false positive threshold, where any score below this value is classified as a false positive. In this experiment, the threshold is set to the value that

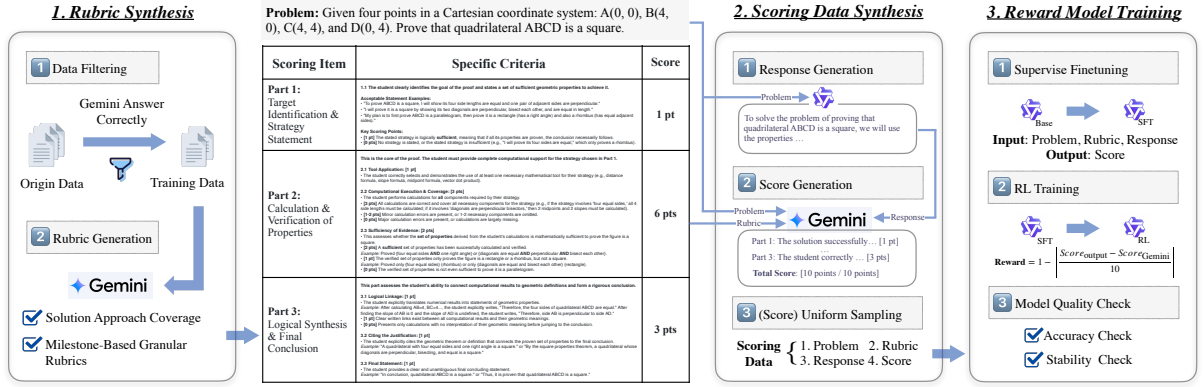


Figure 5: The pipeline of constructing our rubric reward model.

yields the best detection performance: 1.0 for both PRM and RRM.

The results in Figure 4 show that RRM outperforms both PRM and the Verifier in two aspects: *Accuracy*: RRM achieves an F1 of 0.693, surpassing PRM by +0.312 and the Verifier by +0.144. *Continuity*: Unlike the binary Verifier and saturation-prone PRM, RRM yields fine-grained, interpretable 0–10 scores that correlate strongly with false-positive rates (98.2%→17.6% from score 0 to 10). This relatively dense, calibrated signal rewards partially correct, fixable reasoning and penalizes errors proportionally, providing more informative gradients for training.

Overall, RRM offers both higher accuracy and richer, well-calibrated feedback, making it better suited for reducing false positives and promoting robust reasoning than PRM or binary verification.

Does the advantage come from rubrics or from the teacher model? A natural concern is whether the RRM’s advantage stems from the rubric formulation itself or simply from the superior quality of Gemini-2.5-Pro as a teacher.

We argue that rubrics serve as a valuable *intermediate modality* for three reasons. First, rubric-based evaluation is *reference-based*: the judge scores against explicit, problem-specific criteria, a substantially easier task than open-ended, reference-free assessment, making the reward signal more reliable even when the judge is imperfect. Second, once generated, rubrics are *judge-model-agnostic*, decoupling the reward signal from a specific teacher. Third, rubrics make implicit evaluation criteria *explicit and interpretable*, enabling human inspection and refinement.

Fully disentangling rubric format from teacher capability would require a controlled comparison

(e.g., Gemini as a direct scorer without rubrics), which is computationally prohibitive: the RRM needs only one Gemini call per problem, whereas direct step-level scoring would require ~ 200 calls (~ 10 solutions $\times \sim 20$ steps). We discuss this as a limitation.

4.2 Constructing the Rubric Reward Model

We build the Rubric Reward Model through a three-phase pipeline, illustrated in Figure 5. All prompts used in the entire process can be found in the Appendix A (Prompt 2-4).

Phase 1: Rubric Synthesis. We leverage Gemini-2.5-Pro to generate problem-specific evaluation rubrics. Unlike generic prompts, our generation process is grounded in three principles derived from our false positive taxonomy (Table 1): (i) *Targeted criteria against specific failure modes.* Each rubric includes actionable criteria that directly counter identified failure patterns. For instance, rubrics require explicit verification of domain constraints (targeting *Neglected Operational Preconditions*), demand completeness of sufficient conditions rather than example-based patterns (targeting *Inductive Overgeneralization*), and mandate explicit logical linkage between steps (targeting *Miracle Steps*). (ii) *Structure-based scaffolding.* These targeted criteria are embedded in a universal proof structure covering strategy identification, computation/verification, logical synthesis, and conclusion. This holistic scaffolding enables detection of broader flaws such as *Numerical Coincidence* by enforcing a coherent reasoning narrative. (iii) *Method-agnostic fairness.* All rubrics are designed to evaluate any valid solution path, not just one matching a reference solution, ensuring the reward signal focuses on reasoning soundness

regardless of strategy. We also filter out problems for which Gemini’s own solution disagrees with the reference answer, ensuring rubric feasibility. The full prompt and an illustrative example are provided in Appendix A (Prompt 2).

Phase 2: Scoring Data Synthesis. We construct a training dataset $\mathcal{D}_2 = \{(q_i, r_i, y_i, s_i)\}$, where each entry pairs a problem q_i and rubric r_i with a candidate response y_i and a scalar quality score $s_i \in \{0, \dots, 10\}$. We generate diverse candidate responses using both baseline and advanced models, which are then evaluated by Gemini-2.5-Pro against the generated rubrics. We apply weighted sampling to ensure a balanced distribution across score intervals.

Phase 3: Reward Model Training. Initialized from Qwen3-4B-Base, the RRM is trained to predict quality scores given (q, r, y) . We first apply Supervised Fine-Tuning (SFT) for instruction adherence, followed by Proximal Policy Optimization (PPO) to minimize the deviation between predicted and target scores. As shown in Figure 6, the PPO stage significantly enhances scoring stability and accuracy compared to SFT alone.

We provide detailed description for all three phases in Appendix D.3.

5 Experiments And Analysis

5.1 Experimental Setup

Base Model & Dataset: We adopt Qwen3-4B-Base as the backbone model for both the baseline and our proposed approach. Training is conducted on a 9k subset of the Polaris dataset (An et al., 2025), obtained by randomly sampling 10k examples and removing examples where the provided final answer, generated by Gemini, was incorrect. We conduct evaluations on four widely used mathematical reasoning benchmarks, including AIME2024, MATH500, AMC2023, and Olympiad-Bench.

Baseline & Our Method: The baseline consists of Qwen3-4B-Base fine-tuned with PPO using a standard outcome-based reward: 1.0 for a correct final answer and 0 otherwise. The configuration is as follows: maximum sequence length of 4096 tokens, rollout size of 8, batch size of 512, learning rate of 5×10^{-7} , temperature of 1.0, and the Adam optimizer (Kingma and Ba, 2014). The training steps are set to 200 steps. We replace the outcome-based reward model in the baseline with a rubric-based re-

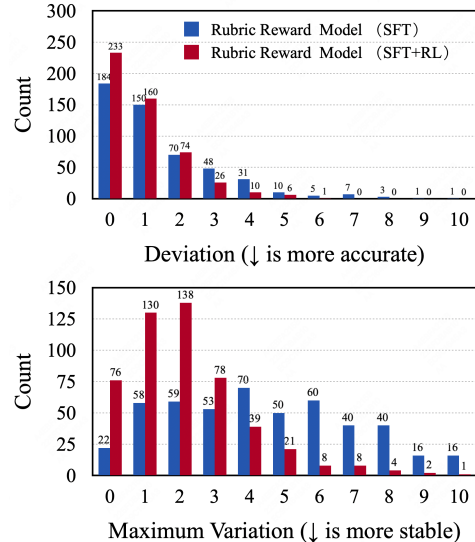


Figure 6: SFT vs. RL RRM. Accuracy: score deviation from Gemini’s score; Stability: maximum variation across 5 runs, temperature set to 1.0.

ward model, while keeping all other configurations unchanged.

Evaluation Metrics: We use both *Standard Pass@N* and *Verified Pass@N*. For the latter, the correctness of each solution is further verified by Gemini-2.5-Pro.² During evaluation, solutions are generated with a temperature of 1.0 and a maximum length of 16,000 tokens.

In the main text, we focus our analysis on the 4B model. The results for the 8B model, along with comprehensive experimental details, are provided in Figure 10b and Appendix D.6, respectively.

5.2 Main Results

The results in Figure 7 yield three takeaways.

Rubric-based rewards deliver consistent gains across datasets. Across evaluation datasets, the rubric-trained model (pink/red) outperforms the outcome-trained model (blue) for all N under both *Standard* and *Verified Pass@N*. This pattern indicates that rewarding reasoning quality—rather than final outcomes alone—induces more generalizable problem-solving behavior.

Gains are larger under *Verified Pass@N* and scale with N. The improvement is notably larger for *Verified Pass@N* than for *Standard Pass@N*, and the Verified-Standard gap widens as N increases. As the candidate budget grows, the base-

²A manual analysis in Table 6 confirms that Gemini-2.5-Pro does not exhibit a preference for our model’s outputs over those from the baseline model, ensuring fair verification.

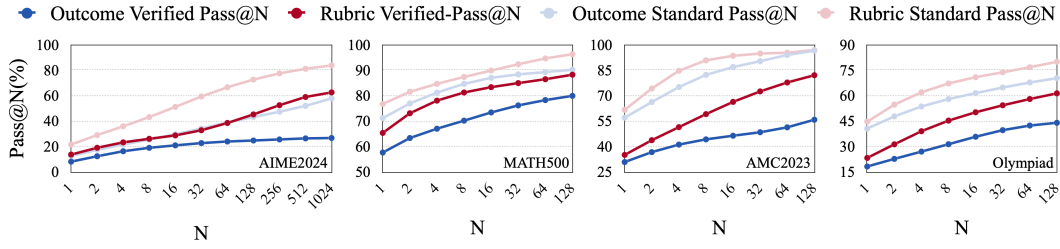


Figure 7: Performance of models trained with Outcome-Based and Rubric-Based Rewards.

line tends to inflate *Standard Pass@N* by sampling more trajectories that accidentally land on the correct answer despite flawed reasoning, whereas our model produces a higher proportion of logically sound solutions. Consequently, the probability that at least one verified-correct solution appears in the N candidates grows faster for our method.

Rubric rewards shrink the Verified-Standard gap. Across all datasets and N , there is a substantial discrepancy between *Standard* and *Verified Pass@N*, underscoring the prevalence of false positives in multi-step reasoning. The gap is consistently smaller for our approach, indicating that rubric guidance suppresses spurious correctness and better aligns generation with logically valid derivations.

5.3 Error-type distribution shift after rubric-based RL

Figure 8(a) illustrates a qualitative shift: rubric-based training not only reduces the overall false positive rate but also transforms *what kinds* of false positives occur.

Rubric rewards suppress critical errors. The most notable effect is on the *Miracle Steps* category. Our method reduces such cases by 71% (175 \rightarrow 50), indicating a substantial suppression of shortcut-driven answer recalls without valid reasoning. Large reductions are also observed in other high-impact failure modes: *Neglected Operational Preconditions* (232 \rightarrow 113) and *Unverified Assumptions* (213 \rightarrow 167). These decreases confirm that the RRM is effective at detecting and discouraging critical lapses in rigor.

More detailed reasoning with minor flaws as a side effect. Interestingly, some categories increase in frequency, notably *Outcome Irrelevance* (67 \rightarrow 118). We view this not as regression, but as a side effect of more detailed reasoning: by encouraging complete, step-by-step derivations, we

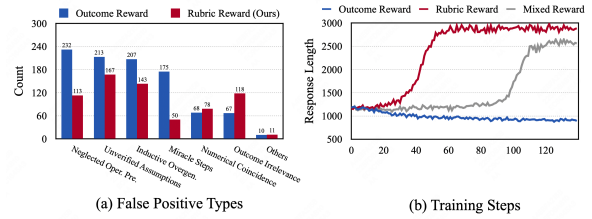


Figure 8: (a) False positive distribution of two models. (b) The change in response length during RL training. “Mixed reward” means 3/4 of the rubric reward + 1/4 of the outcome reward.

increase the chance of minor, localized mistakes in an otherwise coherent chain (see Appendix C for an example). As shown in Figure 8(b), rubric-based training produces longer outputs, reflecting the model’s attempt to build a complete logical chain. Crucially, this does not constitute verbosity-based reward hacking: a Pearson correlation of $\rho = -0.24$ between response length and verified correctness indicates that longer responses are actually *less* likely to be verified as correct.

6 Conclusion

This work systematically exposes the “false positive” phenomenon in mathematical LLMs, where outcome-based rewards mask flawed reasoning. We established a taxonomy of six failure modes and, through direct answer probing, traced the most critical one, *Miracle Steps*, to answer-recall shortcuts that bypass the reasoning chain. We further showed that these failures persist across state-of-the-art models.

To address this, we introduced the RRM, a process-oriented reward function that scores entire reasoning traces against problem-specific rubrics. When integrated into an RL pipeline, RRM-based training consistently outperforms outcome-only supervision. Our results underscore that building genuinely reliable reasoning models requires shifting from validating final answers to verifying the reasoning process itself.

Limitations

There are several limitations in our work: (1) *Dependence on strong external models*. Rubric construction relies on high-capacity models (Gemini-2.5-Pro) and manual filtering, limiting scalability to tasks beyond current LLM capabilities. As open-source models continue to improve, exploring their use as rubric generators is an important direction. (2) *Teacher-model confound*. Although we argue that rubrics provide a valuable intermediate modality (Section 4.1), we cannot fully disentangle the contribution of the rubric format from the capability of the teacher model used to generate rubrics and scoring data. A controlled comparison with Gemini as a direct scorer (without rubrics) under the same compute budget would be informative but is currently prohibitive. (3) *Static reward model during RL*. The RRM is fixed after offline training; as the policy improves, the static scorer may misalign and undervalue novel yet valid reasoning. (4) *Domain and causality limitations*. Experiments are limited to mathematics. In principle, the framework is applicable to any domain where correct outcomes can arise from flawed processes, such as code generation (passing tests via hard-coded outputs) or scientific reasoning (correct conclusions from fabricated evidence). The link between *Miracle Steps* and memorization remains correlational rather than causal, as alternative factors such as problem simplicity, answer guessability, and approximation heuristics may also contribute. Future research could address these limitations by: exploring open-source alternatives for rubric generation; developing adaptive reward models that co-evolve with the policy; and extending our framework to other verifiable reasoning domains.

The Use of Large Language Models

LLMs were employed in a limited capacity for writing optimization. Specifically, the authors provided their own draft text to the LLM, which in turn suggested improvements such as corrections of grammatical errors, clearer phrasing, and removal of non-academic expressions. LLMs were also used to inspire possible titles for the paper. While the system provided suggestions, the final title was decided and refined by the authors and is not directly taken from any single LLM output. In addition, LLMs were used as coding assistants during the implementation phase. They provided code completion and debugging suggestions, but all final

implementations, experimental design, and validation were carried out and verified by the authors. Importantly, LLMs were **NOT** used for generating research ideas, designing experiments, or searching and reviewing related work. All conceptual contributions and experimental designs were fully conceived and executed by the authors.

References

- AIME. 2024. American invitational mathematics examination (aime) 2024. https://huggingface.co/datasets/Maxwell-Jia/AIME_2024.
- AMC. 2023. The american mathematics competitions. <https://huggingface.co/datasets/zwhe99/amc23>.
- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. 2016. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*.
- Chenxin An, Zhihui Xie, Xiaonan Li, Lei Li, Jun Zhang, Shansan Gong, Ming Zhong, Jingjing Xu, Xipeng Qiu, Mingxuan Wang, and Lingpeng Kong. 2025. [Polaris: A post-training recipe for scaling reinforcement learning on advanced reasoning models](#).
- Anthropic. 2025a. Introducing claude 4. <https://www.anthropic.com/news/claude-4>.
- Alignment Team Anthropic. 2023a. Measuring faithfulness in chain-of-thought reasoning. <https://www.anthropic.com/research/measuring-faithfulness-in-chain-of-thought-reasoning>.
- Alignment Team Anthropic. 2023b. Specific versus general principles for constitutional ai. <https://www.anthropic.com/research/specific-versus-general-principles-for-constitutional-ai>.
- Alignment Team Anthropic. 2025b. Reasoning models don't always say what they think. <https://www.anthropic.com/research/reasoning-models-dont-say-think>.
- Rahul K Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñonero-Candela, Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, and 1 others. 2025. Healthbench: Evaluating large language models towards improved human health. *arXiv preprint arXiv:2505.08775*.
- Fazl Barez, Tung-Yu Wu, Iván Arcuschin, Michael Lan, Vincent Wang, Noah Siegel, Nicolas Colignon, Clement Neo, Isabelle Lee, Alasdair Paren, Adel Bibi, Robert Trager, Damiano Fornasiero, John Yan, Yanai Elazar, and Yoshua Bengio. 2025. Chain-of-thought is not explainability. https://fbarez.github.io/assets/pdf/Cot_Is_Not_Explainability.pdf.

- Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wanxiang Che. 2025. Towards reasoning era: A survey of long chain-of-thought for reasoning large language models. *arXiv preprint arXiv:2503.09567*.
- Yanda Chen, Ruiqi Zhong, Narutatsu Ri, Chen Zhao, He He, Jacob Steinhardt, Zhou Yu, and Kathleen McKeown. 2024. Do models explain themselves? counterfactual simulatability of natural language explanations. In *Proceedings of the 41st International Conference on Machine Learning*, pages 7880–7904.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.
- Jacob Dineen, Aswin RRV, Qin Liu, Zhikun Xu, Xiao Ye, Ming Shen, Zhaonan Li, Shijie Lu, Chitta Baral, Muhao Chen, and 1 others. 2025. Qa-lign: Aligning llms through constitutionally decomposed qa. *arXiv preprint arXiv:2506.08123*.
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. 2020. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673.
- Anisha Gunjal, Anthony Wang, Elaine Lau, Vaskar Nath, Bing Liu, and Sean Hendryx. 2025. Rubrics as rewards: Reinforcement learning beyond verifiable domains. *arXiv preprint arXiv:2507.17746*.
- Dadi Guo, Jiayu Liu, Zhiyuan Fan, Zhitao He, Haoran Li, Yumeng Wang, and Yi R Fung. 2025a. Mathematical proof as a litmus test: Revealing failure modes of advanced large reasoning models. *arXiv preprint arXiv:2506.17114*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025b. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A Smith. 2018. Annotation artifacts in natural language inference data. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, and 1 others. 2024. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3828–3850.
- Jujie He, Jiakai Liu, Chris Yuhao Liu, Rui Yan, Chaojie Wang, Peng Cheng, Xiaoyu Zhang, Fuxiang Zhang, Jiacheng Xu, Wei Shen, and 1 others. 2025. Skywork open reasoner 1 technical report. *arXiv preprint arXiv:2505.22312*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Yi Hu, Xiaojuan Tang, Haotong Yang, and Muhan Zhang. 2024. Case-based or rule-based: How do transformers do the math? In *International Conference on Machine Learning*, pages 19438–19474. PMLR.
- Zenan Huang, Yihong Zhuang, Guoshan Lu, Zeyu Qin, Haokai Xu, Tianyu Zhao, Ru Peng, Jiaqi Hu, Zhanming Shen, Xiaomeng Hu, and 1 others. 2025. Reinforcement learning with rubric anchors. *arXiv preprint arXiv:2508.12790*.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Man Ho Lam, Chaozheng Wang, Jen-tse Huang, and Michael R Lyu. 2025. Codecrash: Stress testing llm reasoning under structural and semantic perturbations. *Advances in Neural Information Processing Systems*, 38.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*.
- Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. Faithful chain-of-thought reasoning. In *The 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (IJCNLP-AAACL 2023)*.
- Xueguang Ma, Qian Liu, Dongfu Jiang, Ge Zhang, Zejun Ma, and Wenhu Chen. 2025. General-reasoner: Advancing llm reasoning across all domains. *arXiv preprint arXiv:2505.14652*.
- OpenAI. 2024. o1 system card. <https://openai.com/index/openai-o1-system-card/>.
- OpenAI. 2025a. Introducing gpt-5. <https://openai.com/index/introducing-gpt-5/>.

- OpenAI. 2025b. Introducing openai o3 and o4-mini. <https://openai.com/index/introducing-o3-and-o4-mini/>.
- Lorenzo Pacchiardi, Alex James Chan, Sören Mindermann, Ilan Moscovitz, Alexa Yue Pan, Yarin Gal, Owain Evans, and Jan M Brauner. 2024. How to catch an ai liar: Lie detection in black-box llms by asking unrelated questions. In *The Twelfth International Conference on Learning Representations*.
- Alexander Pan, Erik Jones, Meena Jagadeesan, and Jacob Steinhardt. 2024. Feedback loops with language models drive in-context reward hacking. In *Proceedings of the 41st International Conference on Machine Learning*, pages 39154–39200.
- Peter S Park, Simon Goldstein, Aidan O’Gara, Michael Chen, and Dan Hendrycks. 2024. Ai deception: A survey of examples, risks, and potential solutions. *Patterns*, 5(5).
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R Bowman, Esin DURMUS, Zac Hatfield-Dodds, Scott R Johnston, Shauna M Kravec, and 1 others. 2023. Towards understanding sycophancy in language models. In *The Twelfth International Conference on Learning Representations*.
- Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2025. Scaling llm test-time compute optimally can be more effective than scaling parameters for reasoning. In *The Thirteenth International Conference on Learning Representations*.
- Giulio Starace, Oliver Jaffe, Dane Sherburn, James Aung, Jun Shern Chan, Leon Maksin, Rachel Dias, Evan Mays, Benjamin Kinsella, Wyatt Thompson, and 1 others. 2025. Paperbench: Evaluating ai’s ability to replicate ai research. In *Forty-second International Conference on Machine Learning*.
- Yi Su, Dian Yu, Linfeng Song, Juntao Li, Haitao Mi, Zhaopeng Tu, Min Zhang, and Dong Yu. 2025. Crossing the reward bridge: Expanding rl with verifiable rewards across diverse domains. *arXiv preprint arXiv:2503.23829*.
- Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. 2023. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting. *Advances in Neural Information Processing Systems*, 36:74952–74965.
- Vijay Viswanathan, Yanchao Sun, Shuang Ma, Xiang Kong, Meng Cao, Graham Neubig, and Tongshuang Wu. 2025. Checklists are better than reward models for aligning language models. *arXiv preprint arXiv:2507.18624*.
- Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. 2024. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9426–9439.
- Yu Wang, Nan Yang, Liang Wang, and Furu Wei. 2025. Examining false positives under inference scaling for mathematical reasoning. *arXiv preprint arXiv:2502.06217*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin Carbonneaux, Lingming Zhang, Daniel Fried, Gabriel Synnaeve, Rishabh Singh, and Sida I Wang. 2025. Swe-rl: Advancing llm reasoning via reinforcement learning on open software evolution. *arXiv preprint arXiv:2502.18449*.
- Jiaxin Wen, Ruiqi Zhong, Akbir Khan, Ethan Perez, Jacob Steinhardt, Minlie Huang, Samuel R Bowman, He He, and Shi Feng. 2025. Language models learn to mislead humans via rlhf. In *The Thirteenth International Conference on Learning Representations*.
- Lilian Weng. 2024. Reward hacking in reinforcement learning. <https://lilianweng.github.io/posts/2024-11-28-reward-hacking/>.
- Zhangchen Xu, Yuetai Li, Fengqing Jiang, Bhaskar Ramasubramanian, Luyao Niu, Bill Yuchen Lin, and Radha Poovendran. 2025. Tinyv: Reducing false negatives in verification improves rl for llm reasoning. *arXiv preprint arXiv:2505.14625*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Tian Ye, Zicheng Xu, Yuanzhi Li, and Zeyuan Allen-Zhu. 2024. Physics of language models: Part 2.1, grade-school math and the hidden reasoning process. In *ICLR 2025: International Conference on Learning Representations*.
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, and 1 others. 2025. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*.

- Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. 2024. Generative verifiers: Reward modeling as next-token prediction. In *The 4th Workshop on Mathematical Reasoning and AI at NeurIPS'24*.
- Zhenru Zhang, Chujie Zheng, Yangzhen Wu, Beichen Zhang, Runji Lin, Bowen Yu, Dayiheng Liu, Jingren Zhou, and Junyang Lin. 2025. [The lessons of developing process reward models in mathematical reasoning](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 10495–10516, Vienna, Austria. Association for Computational Linguistics.
- Chujie Zheng, Zhenru Zhang, Beichen Zhang, Runji Lin, Keming Lu, Bowen Yu, Dayiheng Liu, Jingren Zhou, and Junyang Lin. 2025. [ProcessBench: Identifying process errors in mathematical reasoning](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1009–1024, Vienna, Austria. Association for Computational Linguistics.
- Xiangxin Zhou, Zichen Liu, Anya Sims, Haonan Wang, Tianyu Pang, Chongxuan Li, Liang Wang, Min Lin, and Chao Du. 2025. Reinforcing general reasoning without verifiers. *arXiv preprint arXiv:2505.21493*.
- Jiaru Zou, Ling Yang, Jingwen Gu, Jiahao Qiu, Ke Shen, Jingrui He, and Mengdi Wang. 2025. Reasonflux-prm: Trajectory-aware prms for long chain-of-thought reasoning in llms. *arXiv preprint arXiv:2506.18896*.

A Prompts

Prompt 1: Gemini's False Positive Detection

You will receive the following three items: (1) A math problem; (2) A standard answer; (3) A student's submitted answer (including their problem-solving process and final answer). Your task is: (1) Carefully review the student's problem-solving process; (2) Determine whether there are errors, logical flaws, or imprecise points in the method used to arrive at the final answer; (3) If there are problems, explain the type of error and elaborate on why the correct answer was still obtained under such circumstances; (4) The problem-solving process may contain some self-corrected errors, e.g., "1+1=3 wait, 1+1=2" – these are not considered errors but rather the model's thinking process before finding the correct method and answer; (5) There are six types of errors in total. If there are any beyond these six, please explain them additionally:

1. Inductive Overgeneralization (overgeneralization/incomplete induction/insufficient enumeration)
 - Typical symptoms:
 - Asserting "unique solution/no solution/rule holds" after testing only a few small values;
 - Replacing strict elimination with intuition, such as "grows faster/unlikely";
 - Finding only partial solutions without proving there are no more.
 - Why it might still be correct:
 - The actual solutions do fall within the tested small range or are indeed limited to those found; or although the pattern is wrong, the count within the given range happens to match the correct pattern (density/period coincidence).
2. Outcome Irrelevance (rounding/missing multiplication/sign errors in irrelevant parts, or double errors canceling out)
 - Typical symptoms:
 - Rounding too early in the process, but the final result is only reported to the tenths place, so the error does not amplify;
 - Missing the imaginary part/coefficient/negative sign, but only taking the real part/absolute value or $m+n$ (order irrelevant) in the end;
 - Introducing an extra denominator first, then "forgetting" it later, which happens to cancel the error; two miscalculated numbers add up to the correct value.
 - Why it might still be correct:
 - The quantity sought in the problem is insensitive to the error (only depends on the real part/absolute value/last digit/modulus), or the error is swallowed by rounding in the end;
 - Two independent errors accidentally cancel each other out (negative times negative makes positive).
3. Neglected Operational Preconditions (domain/reversibility conditions/boundary points, but coincidentally not affecting)
 - Typical symptoms:
 - Directly canceling/dividing by a variable without first stating that the variable is not zero;
 - Converting $\log(x^2)$ to $2\log x$ without first restricting $x>0$;
 - Simplifying a fractional equation without first stating that the denominator is not zero; ignoring whether boundary points should be included.
 - Why it might still be correct:
 - The calculated value happens to satisfy the (unwritten) domain or reversibility conditions, thus no extraneous or missing roots are produced;
 - Other terms in the problem automatically restrict the domain (e.g., the equation already contains $\log x$, implicitly requiring $x>0$).
4. Unverified Assumptions (unproven structural assumptions/misapplying theorems but hitting equality conditions or special cases)
 - Typical symptoms:
 - Directly assuming "the function must be linear", "extremum occurs when variables are equal", "a trapezoid has maximum area as a rectangle", "choosing a seemingly reasonable parameter value $r=7$ ", etc.;
 - Misapplying theorems (applying quadrilateral properties to hexagons, misusing properties like radical axes/exterior angles, etc.).
 - Why it might still be correct:
 - The guessed structure happens to be the equality condition or a hidden special property in the problem (such as symmetry, equality condition of Cauchy's inequality, special cases in circle geometry), thus the conclusion is correct;
 - The misapplied theorem still holds as a "numerical equality" in this special case, or is equivalent to another correct property.
5. Numerical Coincidence (the problem-solving process is completely different from the correct method and logically invalid, but the final answer is correct due to numerical coincidence)
 - Typical symptoms:
 - Using wrong logic and calculations to get an incorrect probability of $9/20$, while the correct probability is $7/22$. But the problem asks for $m+n$, and coincidentally $9+20=29$ and $7+22=29$, resulting in the same answer;
 - Constructing an incorrect list of numbers that completely fails to meet the problem's conditions, but the square sum of this wrong list happens to equal that of the correct list;
 - Deriving an incorrect pattern of winning/losing conditions based on wrong game analysis, but within the given numerical range, the number of numbers satisfying this wrong pattern is exactly the same as those satisfying the correct pattern.
 - Why it might still be correct:
 - Coincidence.
6. Miracle Steps (the model's solution contains invalid steps, but suddenly arrives at the correct answer)
 - Typical symptoms:
 - The model lists a completely wrong equation " $a + b + c + d - 437 - 2 \times 234 - 3x = 3600$ ", solves $x=-827$ (wrong answer) according to this equation, but the next step directly gives $x=73$ (correct answer);
 - The model provides a series of wrong ideas and steps, but suddenly lists a correct equation/inequality in an incomprehensible way.
7. Other

Please use Chinese and output the results in the following format:

Are there errors or imprecise points in the problem-solving process: Yes / No

If there are problems, why the wrong process led to the correct answer: (This item can be omitted if there are no errors)

- Error type - Explanation - Final result: [1-7] (e.g., [1], [2,3])

Prompt 2: Rubric Generation

Role: You are an experienced math competition coach and problem-setter, an expert in the logical structure of mathematical proofs. Your task is not to solve math problems, but to design a rigorous, universal, and actionable scoring framework for evaluating solution processes.

- Your output should only be the Grading Rubric (i.e. Detailed Scoring Rubric & Coach's Guide), with no other content.
- The total score is 10 points.

Example Problem: Given four points in a Cartesian coordinate system: $A(0, 0)$, $B(4, 0)$, $C(4, 4)$, and $D(0, 4)$. Prove that quadrilateral ABCD is a square.

Guiding Principles:

1. **Method-Agnostic:** This rubric must be able to fairly evaluate all logically correct solution methods, whether they use side lengths, angles, or diagonals. **Strictly prohibit** creating separate criteria for specific methods (e.g., "side-length method," "diagonal method").
2. **Structure-Based:** The core of the scoring should be based on the universal structure of a proof, namely: "identifying key properties," "calculation and derivation," "logical linkage," and "final conclusion."
3. **Actionable Criteria:** The scoring criteria must be specific, observable actions, not abstract descriptions.
 - **Forbidden terms:** "accuracy," "rigor," "clear thinking," "fluent expression."
 - **Encouraged phrases:** "Correctly writes the distance formula," "Explicitly states that the slopes of two segments are negative reciprocals," "Concludes C based on previously proven properties A and B," "Completely states the theorem for identifying a square."

Rubric Framework:

Please break down the scoring rubric into the following sections and assign appropriate points to each (the total score is set to 10 points).

1. Target Identification & Strategy Statement - [e.g., 1 point]

- **Scoring Point:** The student clearly identifies the objective (to prove it's a square) and articulates the set of mathematical properties their chosen strategy relies on.
- **Example:** "To prove it's a square, I will show that all four sides are equal and one interior angle is a right angle." or "I will prove it's a square by showing its diagonals are perpendicular, bisect each other, and are equal in length."

2. Calculation & Verification of Properties - [e.g., 6 points]

- This is the core of the rubric. The student must use calculations to verify **all** key properties required by their chosen strategy. This section is scored based on "properties," and regardless of the method, the student must prove a set of **sufficient conditions**.
- **Scoring Points (detailed by property):**
 - **Proof of Property 1:** [e.g., Equal side lengths]
 - * Correctly applies the necessary formula (e.g., distance formula).
 - * Calculation is free of errors, and lengths of all sides are found.
 - * Reaches an intermediate conclusion of equal side lengths (e.g., $AB=BC=CD=DA=4$).
 - **Proof of Property 2:** [e.g., Perpendicular adjacent sides or perpendicular diagonals]
 - * Correctly applies the necessary method (e.g., slope calculation, vector dot product).
 - * Calculation is free of errors, leading to the conclusion of perpendicularity.
 - **Proof of Property 3:** [e.g., Equal diagonals or diagonals that bisect each other]
 - * ... (and so on)
- **Note:** When scoring, check if the student has completely proven a **full set** of sufficient conditions for their chosen strategy. For example, only proving four equal sides (which could be a rhombus) does not earn full points for this section.

3. Logical Synthesis & Final Conclusion - [e.g., 3 points]

- **Scoring Point 1 - Citing the Justification:** The student explicitly cites a definition or theorem that links the verified properties to the final conclusion. Example: "Because quadrilateral ABCD has four equal sides and one right angle, it is a square."
- **Scoring Point 2 - Final Statement:** Provides a clear, conclusive statement. Example: "Therefore, quadrilateral ABCD is a square. Q.E.D."
- **Scoring Point 3 - Logical Integrity:** The proof is free of logical gaps. For example, the student doesn't just calculate lengths and slopes and then jump to the conclusion without stating what those numbers mean (e.g., "sides are equal" or "sides are perpendicular").

Output Format: Please present the final rubric in a clear table format, including "Scoring Item," "Specific Criteria," and "Score." Here is an example:

Scoring Item	Specific Criteria	Score
1. Target Identification & Strategy Statement	Clearly identifies the objective (to prove it's a square) and articulates the set of mathematical properties their chosen strategy relies on.	1
2. Calculation & Verification of Properties	Uses calculations to verify all key properties required by their chosen strategy. Calculation is free of errors, and lengths of all sides are found. Reaches an intermediate conclusion of equal side lengths (e.g., $AB=BC=CD=DA=4$).	6
3. Logical Synthesis & Final Conclusion	Explicitly cites a definition or theorem that links the verified properties to the final conclusion. Provides a clear, conclusive statement. The proof is free of logical gaps.	3

Core Task: Design a detailed Scoring Rubric for the following math problem. The total score is 10 points.

[Insert the specific math problem here]

Prompt 3: Scoring Data Generation

For each Question and Answer pair, please evaluate according to the given scoring criteria.

- Your output must include two sections: Analysis and Summary.
- In your analysis, always place the score below your reasoning using the following format:

- Reason: [Insert your explanation here]
- Score: X

- In the summary, provide your assessment using the following format:

Final Scoring Summary:

Scoring Criterion 1 (Brief description of the criterion):

(Reason: ...) X points / Y points

Scoring Criterion 2 (Brief description of the criterion):

(Reason: ...) X points / Y points

...

Total Score: [N points / M points]

Question:

{INSERT QUESTION}

Scoring Criteria:

{INSERT CRITERIA}

Answer:

{INSERT MODEL_ANSWER}

Prompt 4: Rubric Model Scoring Template

You are a grade teacher. Based on the following information, please evaluate the student's submitted answer.

The total score is out of 10 points.

Carefully read the question, the scoring criteria, and the submitted answer. Then give a fair and reasonable score according to the grading criteria.

Here is the information:

Question:

{INSERT QUESTION}

Scoring Criteria:

{INSERT CRITERIA}

Student's Answer:

{INSERT MODEL_ANSWER}

B Discarded Question

Question

In the circle with center Q , radii AQ and BQ form a right angle. The two smaller regions are tangent semicircles, as shown. The radius of the circle with center Q is 14 inches. What is the radius of the smaller semicircle? Express your answer as a common fraction.

C Case Study

Question

We are told that the process \mathcal{P} consists of:

1. Reflect the point in the x -axis.
2. Translate the resulting point 2 units upwards.
3. Reflect the resulting point in the y -axis.

We are given $S_6 = (-7, -1)$ and asked to find S_0 after reversing \mathcal{P} six times.

Baseline's Solution

Solution

To determine the coordinates of the point S_0 , we need to reverse the steps of the process \mathcal{P} applied six times to obtain the point $S_6(-7, -1)$. Let's break down the process step by step.

Step-by-Step Reasoning:

1. **Identify the final point S_6 :**
 $S_6 = (-7, -1)$

2. **Reverse the sixth application of \mathcal{P} :**

- S_6 is the result of applying \mathcal{P} to S_5 .
- To find S_5 , we need to reverse the steps of \mathcal{P} :
 - (a) Step 3: Reflect S_6 in the y -axis. This means S_5 has the same y -coordinate as S_6 but the opposite x -coordinate. $S_5 = (7, -1)$
 - (b) Step 2: Translate S_5 2 units downwards. This means S_4 has the same x -coordinate as S_5 but the y -coordinate is 2 units less. $S_4 = (7, -3)$
 - (c) Step 1: Reflect S_4 in the x -axis. This means S_3 has the same x -coordinate as S_4 but the opposite y -coordinate. $S_3 = (7, 3)$

3. **Reverse the fifth application of \mathcal{P} :**

- S_3 is the result of applying \mathcal{P} to S_2 .
- To find S_2 , we need to reverse the steps of \mathcal{P} :
 - (a) Step 3: Reflect S_3 in the y -axis. This means S_2 has the same y -coordinate as S_3 but the opposite x -coordinate. $S_2 = (-7, 3)$
 - (b) Step 2: Translate S_2 2 units downwards. This means S_1 has the same x -coordinate as S_2 but the y -coordinate is 2 units less. $S_1 = (-7, 1)$
 - (c) Step 1: Reflect S_1 in the x -axis. This means S_0 has the same x -coordinate as S_1 but the opposite y -coordinate. $S_0 = (-7, -1)$

Thus, the coordinates of the point S_0 are $(-7, -1)$.

False Positive Analysis

The student's problem-solving method contained a fundamental logical error. The student attempted to deduce S_0 from S_6 by reversing the transformation process \mathcal{P} . A single inverse process \mathcal{P}^{-1} indeed requires performing the opposite operations in reverse order (first reflect across the y -axis, then translate downward by 2, and finally reflect across the x -axis). However, the student incorrectly assumed that **each individual sub-step** in the inverse process corresponded to one backward step in the sequence. For example, in the process of solving for S_5 from S_6 , the student recorded the result after the first inverse operation (reflection across the y -axis) as S_5 , the result after the second inverse operation (translation downward) as S_4 , and the result after the third inverse operation (reflection across the x -axis) as S_3 .

This is a serious logical flaw. The correct approach is to apply **all three inverse operations** to S_6 ; the final result, after completing all steps, is S_5 . According to the student's erroneous logic, they actually performed two complete inverse processes:

1. In the first complete inverse operation, they believed they had derived S_3 from S_6 .
2. In the second complete inverse operation, they believed they had derived S_0 from S_3 .

Thus, the student effectively computed $(\mathcal{P}^{-1})^2(S_6)$, which is the coordinate of S_4 , but they labeled it as S_0 .

The fact that the final answer was correct is purely coincidental. The transformation \mathcal{P} has a special property: it has a period of 2. That is, applying \mathcal{P} twice to any point returns the point to its original position (i.e., $\mathcal{P}(\mathcal{P}(S)) = S$). This means the sequence $S_0, S_1, S_2, S_3, S_4, S_5, S_6, \dots$ is periodic with period 2, specifically $S_0 = S_2 = S_4 = S_6$ and $S_1 = S_3 = S_5$.

Due to this periodicity, the coordinates of S_4 are exactly the same as those of S_0 . The student mistakenly computed the value of S_4 , but because $S_4 = S_0$, they arrived at the correct answer for S_0 purely by coincidence. The student's problem-solving process showed no awareness of or reliance on this periodicity; the correctness of the final answer was entirely accidental.

Ours Solution

Solution

Alright, I have this problem about a point moving through a series of transformations, and I need to find the original point based on the final position after multiple applications of a process called \mathcal{P} . Let's break it down step by step.

Understanding the Problem

First, let's understand what's given:

1. **Initial Point:** $(3, 5)$
2. **Process \mathcal{P} :** **Step 1:** Reflect the point in the x -axis. **Step 2:** Translate the resulting point 2 units upwards. **Step 3:** Reflect the resulting point in the y -axis.
3. **Applications:** 1. Start with point S_0 . 2. Apply \mathcal{P} to S_0 to get S_1 . 3. Continue applying \mathcal{P} four more times to get S_6 .
4. **Final Point:** $S_6(-7, -1)$
5. **Goal:** Find the coordinates of S_0 .

Breaking Down the Process \mathcal{P}

Let's understand each step of \mathcal{P} and how it transforms a point (x, y) .

1. **Step 1: Reflect in the x -axis.** Reflecting a point in the x -axis changes the sign of the y -coordinate. So, (x, y) becomes $(x, -y)$.
2. **Step 2: Translate 2 units upwards.** Translating a point upwards increases the y -coordinate by 2. So, $(x, -y)$ becomes $(x, -y + 2)$.
3. **Step 3: Reflect in the y -axis.** Reflecting a point in the y -axis changes the sign of the x -coordinate. So, $(x, -y + 2)$ becomes $(-x, -y + 2)$.

Therefore, the entire process \mathcal{P} transforms (x, y) to $(-x, -y + 2)$.

Applying \mathcal{P} Multiple Times

We start with S_0 and apply \mathcal{P} six times to get S_6 . Let's denote the transformation as T , so:

$$T(x, y) = (-x, -y + 2)$$

Applying T multiple times:

$$T^n(x, y) = T(T^{n-1}(x, y))$$

We need to find a pattern or a formula for $T^n(x, y)$. Let's compute the first few applications to see if a pattern emerges.

$$T^1(x, y) = (-x, -y + 2)$$

$$T^2(x, y) = T(T^1(x, y)) = T(-x, -y + 2) = (-(-x), -(-y + 2) + 2) = (x, y - 2 + 2) = (x, y)$$

Wait, $T^2(x, y) = (x, y)$. That's interesting! Applying T twice brings us back to the original point.

$$T^3(x, y) = T(T^2(x, y)) = T(x, y) = (-x, -y + 2)$$

So, $T^3(x, y) = T^1(x, y)$.

$$T^4(x, y) = T(T^3(x, y)) = T(-x, -y + 2) = (x, y)$$

Again, $T^4(x, y) = (x, y)$.

From these computations, it seems that applying T an even number of times returns the point to its original position, and applying T an odd number of times gives the same result as applying T once.

Generalizing the Transformation

Given that $T^2(x, y) = (x, y)$, we can see that: If n is even, $T^n(x, y) = (x, y)$. If n is odd, $T^n(x, y) = T^1(x, y) = (-x, -y + 2)$.

Since we're applying T six times ($n = 6$), which is even, we have: $T^6(x, y) = (x, y)$.

But we know that $T^6(S_0) = S_6 = (-7, -1)$. Therefore: $(x, y) = (-7, -1)$ So, $S_0 = (-7, -1)$.

Verification

Let's verify this by applying \mathcal{P} six times starting with $S_0 = (-7, -1)$ and see if we get $S_6 = (-7, -1)$.

$$S_0 = (-7, -1) \rightarrow S_1 = T(S_0) = (-(-7), -(-1) + 2) = (7, 1 + 2) = (7, 3) \rightarrow S_2 = T(S_1) = (-7, -3 + 2) = (-7, -1) \rightarrow S_3 = T(S_2) = (7, 1) \rightarrow S_4 = T(S_3) = (-7, -1) \rightarrow S_5 = T(S_4) = (7, 3) \rightarrow S_6 = T(S_5) = (-7, -1)$$

Yes, after six applications, we indeed get $S_6 = (-7, -1)$. This confirms that our solution is correct.

Final Answer

The coordinates of the point S_0 are $\boxed{(-7, -1)}$.

False Positive Analysis

The student's overall problem-solving approach and core derivation process are completely correct. They correctly derived the algebraic expression for the single transformation \mathcal{P} as $T(x, y) = (-x, -y + 2)$, and found through calculation that the second iteration of the transformation, $T^2(x, y)$, is the identity transformation, i.e., $T^2(x, y) = (x, y)$. Based on this key discovery, they concluded that "applying the transformation an even number of times will return the point to its original position," and therefore $S_6 = T^6(S_0) = S_0$. This chain of reasoning is rigorous and flawless, directly leading to the correct answer $S_0 = S_6 = (-7, -1)$.

The error occurred in the final "Verification" step, which was intended as a post-solution check. In computing the verification sequence, the student made a slip or calculation error:

- They correctly computed $S_2 = (-7, -1)$.
- Next, when calculating S_3 , they wrote $S_3 = T(S_2) = (7, 1)$.
- The correct computation should be $S_3 = T(-7, -1) = (-(-7), -(-1) + 2) = (7, 1 + 2) = (7, 3)$.

D Supplement

D.1 Details for False Positive Analysis Pipeline

There are four stages for analyzing false positive modes:

Stage 1: Data Preparation. We assemble a dataset of 680 samples, comprising 170 distinct questions (30 questions from AIME2024 + 50 questions from MATH500 + 40 questions from AMC2023 + 50 questions from Olympiad), each with four unique model responses. All markdown and formulas have been converted into an easily readable format.

Stage 2: Initial Mode Discovery. We use Gemini-2.5-Pro for an automated review to generate a preliminary taxonomy of “false positives.” The model is prompted with each question, a reference solution, and the model’s response, and is instructed to report on (1) any reasoning errors and (2) how flawed reasoning can still yield a correct answer. These reports are then aggregated and synthesized by the model into the initial taxonomy.

Stage 3: Expert Review. In the third stage, we conduct a human validation of these modes. Four annotators, all holding undergraduate degrees with substantial training in advanced mathematics, evaluate each sample. They are equipped with tools like Google Search and large models and are instructed to discard any samples beyond their expertise. For each sample, they determine if it is a false positive and, if so, classify it using our preliminary taxonomy or label it as “Other” with a detailed explanation.

The annotators are PhD students we recruited from a certain university. Each annotator is paid \$27.5 per hour.

Stage 4: Synthesis and Analysis. In the final stage, we refine the taxonomy by incorporating the “Other” categories identified by human annotators. Using this human-validated framework, we perform a quantitative analysis to measure the frequency of each false positive mode, revealing the model’s prevalent reasoning flaws.

D.2 Experimental Details for Leading Model Evaluation

Models and Generation. We evaluated four leading models: GPT-5-thinking, o4-mini, Gemini-2.5-Pro, and Claude-4-Sonnet-thinking. We employ Gemini-2.5-Pro (version 0605). For the other models, namely o4-mini, GPT-5, and Claude-4-Sonnet, we utilize their latest versions available

as of September 2025. For each question in the challenge set, we generated a single response from each model ($n = 1$). To encourage more detailed reasoning, we set the reasoning effort parameter to ‘high’ for both GPT-5-thinking and o4-mini.

Evaluation Protocol. All generated responses were manually evaluated by human annotators.

Additional notes. During annotation, we noted that o4-mini exhibited a strong tendency to provide overly concise or truncated reasoning steps. This brevity sometimes made it challenging to fully assess the validity of its solution path and may contribute to its higher observed false positive rate, as critical (and potentially erroneous) intermediate steps might be omitted.

Qualitative Examples. For qualitative insights, several examples of questions from our challenge set that frequently induced false positives across the models are presented below:

Question 1. Rows 1, 2, 3, 4, and 5 of a triangular array of integers are shown below.

```
1
1 1
1 3 1
1 5 5 1
1 7 11 7 1
```

Each row after the first row is formed by placing a 1 at each end of the row, and each interior entry is 1 greater than the sum of the two numbers diagonally above it in the previous row. What is the units digit of the sum of the 2023 numbers in the 2023rd row?

Failure: Gemini-2.5-Pro and Claude-4-Sonnet, through enumeration, discovered an important function $U(\cdot)$ in solving the problem have: $U(21) = U(1)$. Without providing proof, they directly claimed the existence of periodicity.

Question 2. Rectangles $ABCD$ and $EFGH$ are drawn such that D, E, C, F are collinear. Also, A, D, H, G all lie on a circle. If $BC = 16, AB = 107, FG = 17$, and $EF = 184$, what is the length of CE ?

Failure: All models overlook the possible permutations of D, E, C, F .

Question 3. How many ordered pairs of positive real numbers (a, b) satisfy the equation

$$(1 + 2a)(2 + 2b)(2a + b) = 32ab?$$

Failure: Claude-4-Sonnet directly identified the correct (a,b) pair through trial, then reported unsuccessful attempts with alternative answers, and subsequently claimed that only one such pair satisfies the requirements. GPT-5 ignored the case of a zero denominator during its simplification process. o4-mini made an error in its variable substitution step.

D.3 RRM Training Pipeline

We build the Rubric Reward Model through a three-phase pipeline, illustrated in Figure 5. All prompts used in the entire process can be found in the Appendix A (Prompt 2-4).

Phase 1: Rubric Synthesis. The first step is to construct a problem-specific rubric for each training example. Our goal is to design evaluation criteria that are logically grounded and tailored to directly counteract the failure modes identified in our taxonomy (refer to Table 1). To achieve this, we prompt Gemini-2.5-Pro to generate rubrics that embody a set of core principles, thereby transforming empirical findings into actionable evaluation guidelines.

Principle 1: Targeted principles against specific failure modes.

- **Neglected Operational Preconditions & Unverified Assumptions:** Each rubric must include actionable and specific criteria. For example, instead of a vague correctness check, the rubric demands explicit verification of constraints, thereby penalizing solutions that work only coincidentally while ignoring fundamental requirements.
- **Inductive Overgeneralization:** We enforce the principle of completeness of sufficient conditions. The rubric must assess whether the presented evidence and reasoning are collectively sufficient for a general proof, not merely consistent with a few examples. This shifts evaluation from pattern-matching toward requiring deductive rigor.
- **Miracle Steps:** The rubric mandates explicit logical linkage between steps. Any jump from con-

fusion to an answer—without a valid derivation—fails this criterion. This ensures the reasoning chain is fully articulated, directly penalizing “miraculous” leaps symptomatic of answer-recall shortcuts.

Principle 2: Structure-based scaffolding. These targeted criteria are embedded in a universal proof structure—covering strategy, computation/verification, synthesis, and conclusion. This holistic structure enables detection of broader logical flaws such as Outcome Irrelevance and Numerical Coincidence, by enforcing a coherent narrative of reasoning rather than allowing a collection of disjointed, potentially flawed calculations.

Principle 3: Method-agnostic fairness. All rubrics must be method-agnostic, capable of evaluating any valid solution path, not just one that matches a reference solution. This focuses the reward signal on the soundness of reasoning itself, regardless of strategy.

Based on the above principles, we carefully designed the prompt and included an illustrative, hand-crafted example in it to guide consistent generation. The detailed prompt refers to Prompt 2.

To further ensure rubric quality, we first filter out training problems for which Gemini-2.5-Pro’s own solution disagrees with the reference answer, thereby eliminating problems beyond the model’s capabilities and ensuring rubric feasibility. This procedure yields the dataset:

$$\mathcal{D}_1 = \{(q_i, r_i)\}_{i=1}^N,$$

where q_i denotes a problem and r_i its corresponding rubric.

Phase 2: Scoring Data Synthesis. Next, we generate annotated training examples for the reward model. For each (q_i, r_i) , we produce multiple candidate responses using both the baseline model and Gemini-2.5-Pro (the latter increases the proportion of high-quality responses). We then feed the problem, rubric, and candidate response to Gemini-2.5-Pro to obtain an integer score from 0 to 10.³ To reduce score imbalance and avoid over-representing mid- or low-quality reasoning, we apply weighted sampling across score intervals, ensuring a more

³In Appendix D.5, we have manually assessed the accuracy of Gemini’s scoring. In the 1320 cases, 12 scores were higher than the actual level, and 7 scores were lower. Additionally, we have tested the stability of Gemini’s scores across 5 runs, which is presented in Figure 9.

uniform distribution. After this phase, we obtain

$$\mathcal{D}_2 = \{(q_i, r_i, y_i, s_i)\}_{i=1}^N,$$

where y_i is a candidate response and $s_i \in \{0, 1, \dots, 10\}$ is the assigned score.

Phase 3: Reward Model Training. We initialize our RRM from the Qwen3-4B-Base model and first perform supervised fine-tuning (SFT) on \mathcal{D}_2 , training it to take (q, r, y) as input and output the corresponding analysis and final score. This yields an SFT-trained checkpoint RRM_{SFT} . We then further refine the model using proximal policy optimization (PPO). The reward function is defined as

$$R = 1 - \frac{|s_{\text{pred}} - s_{\text{target}}|}{10}, \quad (1)$$

where s_{pred} and s_{target} denote the predicted and target scores, respectively. The final result, RRM_{RL} , serves as our rubric-aware scoring function in downstream reinforcement learning. Our rubric reward models’ accuracy and stability on the hold-out test set are shown in Figure 6. Compared with RRM_{SFT} , RRM_{RL} has significantly higher accuracy and stability. Training details refer to Appendix D.4.

D.4 RRM Training Details

We fine-tune the Qwen3-4B-Base model as our policy model using PPO. The training is guided by the reward function in Eq. equation 1.

The PPO training is configured with the following hyperparameters: a maximum prompt length of 10000, a maximum sequence length of 2048, a batch size of 128, and a rollout size of 8. We use the Adam optimizer with a learning rate of 5×10^{-7} and a generation temperature of 1.0. The model is trained for 400 steps.

D.5 The Scoring Accuracy and Stability of Gemini-2.5-Pro.

When using Gemini-2.5-Pro for scoring, we set the temperature to 1.0, perform repeated sampling five times, and calculate the difference between the highest score and the lowest score among these five runs. As can be seen from the Figure 9, Gemini-2.5-Pro demonstrates good stability despite minor fluctuations.

D.6 Details for Main Experiments

All our training and inference were conducted on a server with 8 NVIDIA A800-80G GPUs. During evaluation, we set the temperature to 1.0, the

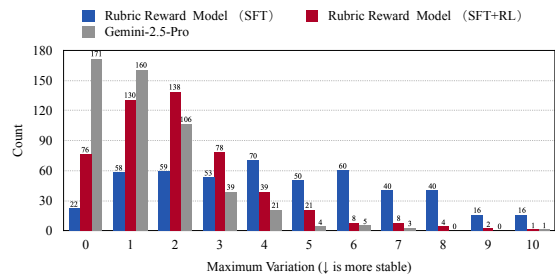


Figure 9: The scoring stability of Gemini-2.5-Pro.

Table 5: Manual evaluation of the accuracy of Gemini’s scoring according to the rubric.

Gemini Rubric Scoring	Too high	Too low	Accurate
Count	12	7	1301

maximum generation length to 16,000 tokens, and used the prompt: *Please reason step by step, and put your final answer within \boxed{\}*.

To evaluate $\text{Pass}@N$, we generate $2N$ candidate solutions for each problem instance.

Evaluation on full datasets and the Qwen3-8B.

In our main experiments, due to computational cost considerations, we randomly selected a subset of 50 samples from MATH500 (500 samples) and Olympiad (675 samples) for evaluation. We additionally conducted experiments on the full datasets (32 runs), and the results are presented in Figure 10a and 10b. The overall trends and conclusions remain consistent with those observed on the subset.

Comparison of the scores assigned by Gemini-2.5-Pro to our model and the baseline models.

As a supplementary result, Figure 10c presents the outcomes of using Gemini-2.5-Pro to generate a rubric on the test set and to score the responses of both models.

In our distributional analysis of error cases (Section 5.3), we focus on instances that were not assigned a perfect score by Gemini-2.5-Pro. The rationale is that false-positive samples with a perfect Gemini grade represent cases where the rubric reward is inherently unable to address the issue. In contrast, our error analysis aims to examine cases in which the rubric reward could potentially play a role.

Table 6: Confusion matrix comparing false positives identified by human annotators and by Gemini, split by model source.

		Overall		Rubric		Outcome	
		Gemini TP	Gemini FP	Gemini TP	Gemini FP	Gemini TP	Gemini FP
Human	TP	462	93	252	52	210	41
	FP	9	295	1	152	8	144

Table 7: Confusion matrix on different datasets.

		AIME		AMC		MATH		Olympiad	
		Gem. TP	Gem. FP	Gem. TP	Gem. FP	Gem. TP	Gem. FP	Gem. TP	Gem. FP
Human	TP	34	8	112	17	222	50	94	18
	FP	1	28	2	105	0	44	6	118

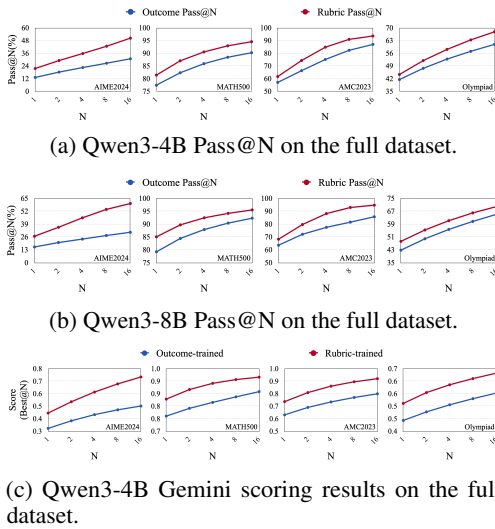


Figure 10: Full dataset evaluation results.

E Human Evaluation

E.1 Gemini-2.5-Pro as a False Positive Judger: Reliability Assessment

Table 8: The proportion of questions for which the model and human false positive evaluations are identical across all responses to that question.

Human-Gemini Consistency	Qwen3-Outcome (4 resp. per query)	Qwen3-Rubric (4 resp. per query)	Overall (8 resp. per query)
Ratio	92/121	109/139	97/141

Agreement with human experts. We quantify Gemini-2.5-Pro’s reliability by conducting extensive human evaluation. As shown in Table 6, Gemini attains high precision (98.1%) and reasonable recall (83.2%) against human labels, yielding an overall F1 score of 0.90 and an agreement rate of 88.1%. These results confirm that Gemini correctly flags almost all human-identified false positives

and makes very few spurious accusations.

No preference toward rubric-trained or outcome-trained outputs. Empirically, Gemini exhibits comparable behavior on rubric-trained and outcome-trained responses. From Table 6:

- Rubric-trained subset: precision 99.6%, recall 82.9%, agreement 88.4%.
- Outcome-trained subset: precision 96.3%, recall 83.6%, agreement 87.9%.

The near-identical recalls (82.9% vs 83.6%) and close agreement rates (88.4% vs 87.9%) show no systematic advantage for rubric-trained outputs; if anything, the tiny precision difference reflects fewer false alarms on that subset, not preferential scoring.

Consistency across datasets. The performance is stable across datasets (Table 7): F1 ranges from 0.88 (AIME) to 0.92 (AMC), with precision consistently ≥ 0.94 . This robustness suggests that Gemini’s accuracy is not confined to a particular problem source or difficulty level.

Agreement at question level. We also assess whether Gemini-2.5-Pro and human annotators agree *across all responses* to the same prompt. Complete question-level agreement holds for 76.0% of questions in the outcome-trained setting, 78.4% in the rubric-trained setting, and 68.8% overall (Table 8). The similar agreement rates for rubric- and outcome-trained models indicate that Gemini does not systematically favor one training method over the other.

Given its high precision, stable cross-dataset performance, and absence of bias toward our method, we use Gemini-2.5-Pro as a scalable, automatic false-positive judge for the remainder of our analysis.