

PRMBENCH: A Fine-grained and Challenging Benchmark for Process-Level Reward Models

Mingyang Song^{1,2}, Zhaochen Su³, Xiaoye Qu², Jiawei Zhou^{4†}, Yu Cheng^{5†}

¹Fudan University, ²Shanghai AI Laboratory, ³Soochow University

⁴Stony Brook University, ⁵The Chinese University of Hong Kong

mysong23@m.fudan.edu.cn; suzhaochen0110@gmail.com;

quxiaoye@pjlab.org.cn; jzhou@ttic.edu; chengyu@cse.cuhk.edu.hk;

Project Page: <https://prmbench.github.io>

Abstract

Process-level Reward Models (PRMs) are crucial for complex reasoning and decision-making tasks, where each intermediate step plays an important role in the reasoning process. Since large language models (LLMs) suffer from various types of errors during the reasoning process, PRMs are required to possess nuanced capabilities for detecting various implicit error types in real-world scenarios. However, current benchmarks primarily focus on step correctness, failing to evaluate PRMs' performance systematically. To address this gap, we introduce PRMBENCH, a process-level benchmark specifically designed to assess the fine-grained error detection capabilities of PRMs. PRMBENCH comprises 6,216 carefully designed problems and 83,456 step-level labels, evaluating models across multiple dimensions, including *simplicity*, *soundness*, and *sensitivity*. In our experiments on 25 models, spanning across both open-source PRMs and LLMs prompted as critic models, we uncover significant weaknesses in current PRMs. These findings reveal the challenges inherent in process-level evaluation and highlight key directions for future research, establishing PRMBENCH as a robust testbed for advancing research on PRM evaluation and development.

1 Introduction

Recent large language models (LLMs) (OpenAI, 2024a,b; Team, 2024a), trained on large-scale reinforcement learning, have achieved significant performance in complex reasoning tasks such as mathematics and code generation (Yu et al., 2023; Guo et al., 2024; DeepMind, 2024; Luo et al., 2023; Qu et al., 2025). A key factor behind their successes is the use of process reward models (PRMs) (Wang et al., 2023; Lightman et al., 2023; Uesato et al., 2022), which can help evaluate the correctness of

[†]Equal senior contribution

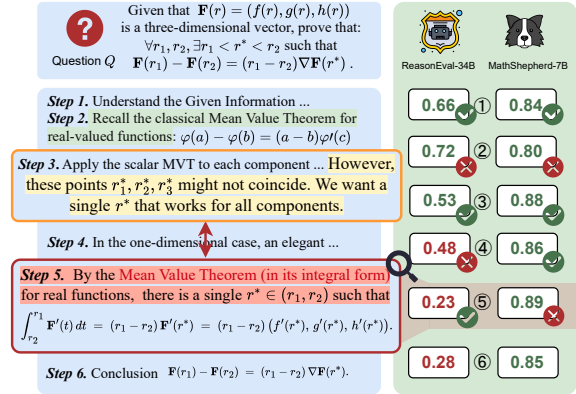


Figure 1: (Left): Given a question Q , the reasoning step 2 and 5 of OpenAI-o1 model contains errors. (Right): The step-level reward scores generated by ReasonEval-34B (Xia et al., 2024b) and MathShepherd-7B (Wang et al., 2023). Green scores indicate the PRM prefer labeling this step as correct while red scores indicate the PRM prefer identifying this step as incorrect.

reasoning steps and train LLMs with appropriate rewards. (Qin et al., 2024; Zhang et al., 2024b).

However, during the reasoning process of the LLM, it suffers from various types of errors, while recent PRMs are not able to identify all these error types precisely. For instance, as illustrated in Figure 1, given a question Q , OpenAI o1 model (OpenAI, 2024b) generates a reasoning procedure containing errors, where step 2 is redundant, step 5 is inconsistent with step 3, and the theory used in step 5 is incorrect due to deception falling. Under these circumstances, ReasonEval-34B (Xia et al., 2024b) and Math-Shepherd-7B (Wang et al., 2023) fail to identify these errors accurately. Math-Shepherd-7B fails to recognize step 5 as an error, while ReasonEval-34B correctly identifies step 5 but incorrectly classifies step 4 as an error, indicating the unreliability of current PRMs.

To evaluate the diverse error-detection capabilities of PRMs, we present PRMBENCH, a comprehensive and fine-grained benchmark specifically

	PRM Benchmarks?	Error Type Detection?	Fine-grained classes [†]	Step Evaluation	Annotator	Test Case Size	Average Steps
MR-GSM8K (Zeng et al., 2023)	✗	✗	1	✓	Human	2,999	8.3
RMBench (Liu et al., 2024)	✗	✗	1	✗	Synthetic + Human	1,327	-
CriticBench (Lin et al., 2024)	✗	✗	1	✗	-	-	-
MathCheck-GSM (Zhou et al., 2024)	✗	✗	1	✓	Synthetic	516	-
MR-Ben (Zeng et al., 2024b)	✗	✗	1	✓	Human	5,975	9.5
ProcessBench (Zheng et al., 2024)	✓	✗	1	✓	Human	3,400	7.1
PRMBENCH	✓	✓	9	✓	Synthetic + Human	6,216	13.4

Table 1: Comparison between our proposed PRMBENCH and other benchmarks or datasets related to reasoning process assessment. [†]: Fine-grained classes mean the number of evaluation categories according to fine-grained error types of model generation.

designed for assessing PRMs. In contrast to existing process-level benchmarks, which can only evaluate the detection of a single error type (Zheng et al., 2024; Zeng et al., 2024b), PRMBENCH offers a more nuanced evaluation. Specifically, PRMBENCH systematically assesses the performance of PRMs across diverse error categories, including *simplicity*, *soundness*, and *sensitivity*. Our benchmark includes 6,216 fine-grained data instances spreading across three major evaluation categories and nine sub-categories, whose quality is ensured by professional annotators. Additionally, we utilize style-controlled data curation methods to ensure evaluation samples under consistent difficulty levels, mitigating confounding variables.

In our study, we conduct extensive experiments using PRMBENCH to evaluate 25 models, including dedicated PRMs and SOTA general-purpose or mathematical LLMs, prompted as critic models. We observe that all PRMs partially grasp multi-step process evaluation. Specifically, Gemini-2-Thinking achieves the best performance of 68.8, but still significantly falls behind the human performance of 83.8. Through extensive analysis, we discover a significant inconsistency between step-level and outcome-level evaluations. By evaluating models with PRMBENCH, we can assess PRMs’ ability to detect step-level errors and false positives, reducing the risk of outcome hacking (Gao et al., 2025). To sum up, our contributions are as follows:

- We present PRMBENCH, the first comprehensive process-level reward model benchmark, comprising 6,216 carefully curated samples and 83,456 step-level labels for a series of evaluations on process-level reward models.
- PRMBENCH covers three carefully-crafted evaluation categories and nine sub-categories including *simplicity*, *soundness*, and *sensitivity*. With these fine-grained evaluation axes, we can conduct tailored assessments of models on their spec-

ified capabilities and reveal their potential weaknesses during the rewarding procedure.

- Based on our proposed PRMBENCH, we conduct in-depth pilot experiments on twenty-five models including PRMs along with SOTA LLMs. Our findings uncover critical weaknesses and provide valuable insights to guide future research to improve the capabilities of PRMs.
- To facilitate future research, we release the PRM-EVAL toolkit, offering an automated evaluation framework and customizable data generation system. We hope PRMBENCH will drive progress in step-level reasoning for RLHF and foster further development of more reliable PRMs.

2 Related Work

2.1 Process-level Reward Models

Process-level reward models (PRMs) have shown improvements over traditional outcome-level reward models (ORMs) in enhancing process-level reasoning accuracy and long-process reasoning abilities (Lightman et al., 2023; Uesato et al., 2022). Recently, several PRMs have been proposed for process-level RLHF (Wang et al., 2023; Xia et al., 2024b; o1 Team, 2024), with Lightman et al. (2023) releasing a large dataset for multi-step reasoning, and Wang et al. (2023) introducing an automatic self-supervised pipeline for process-level labeling. Xia et al. (2024b) uses PRMs as auto-evaluators for multi-step reasoning accuracy. As PRM training and data curation have grown, numerous PRMs (o1 Team, 2024; Xiong et al., 2024; Team, 2024b; Gao et al., 2024) have emerged, along with critic models using LLM-generated feedback (McAleese et al., 2024; Zhang et al., 2024a; Gao et al., 2024). However, both PRMs and critic models remain fallible, highlighting the need for comprehensive benchmarks. In this paper, we propose PRMBENCH, a comprehensive benchmark for evaluating PRMs on

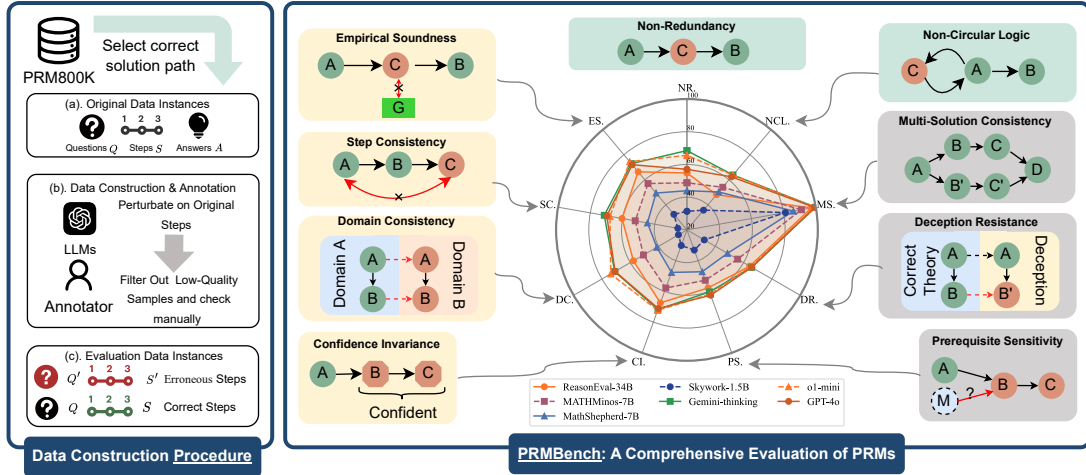


Figure 2: An overview of our PRMBENCH. The left part illustrates our data curation procedure. In the right part of the figure, we showcase demonstrations of our evaluation categories and the relative performance of tested models, with green, yellow, and gray boxes indicating *simplicity*, *soundness*, and *sensitivity* respectively, where red circles represent erroneous steps and green circles indicate correct regular steps.

fine-grained subjects, establishing a strong foundation for PRM evaluation.

2.2 Reasoning Benchmarks

Evaluating the reasoning capabilities of LLMs is crucial for understanding their potential and limitations. ROSCOE (Golovneva et al., 2022) introduces a semantic comparison-based multi-step reasoning accuracy evaluation benchmark. However, recent research suggests that labeled data cannot be assumed to cover all possible solution paths exhaustively (Wang et al., 2023; Xia et al., 2024b). To address this, Xia et al. (2024b) uses PRMs or Critic models to evaluate step-level reasoning accuracy. However, PRMs are not always accurate in assessing process-level data, underscoring the need for a comprehensive evaluation benchmark. While other benchmarks (Liu et al., 2024; Li et al., 2024a; Lin et al., 2024; Su et al., 2024b) exist, they are not tailored for PRMs and can’t assess step-level reasoning. Some works (Zeng et al., 2023, 2024b; Yan et al., 2024) use LLMs to evaluate reasoning steps, but they often overlook implicit error types. Existing error classification works are not specific to PRMs and lack fine-grained step-level labels (Li et al., 2024b). To address these gaps, we propose PRMBENCH, a solution that offers fine-grained evaluation and detects various error types.

3 PRMBENCH

3.1 Evaluation Categories

In this section, we provide a detailed introduction to the evaluation categories of PRMBENCH, which

is organized into three main domains:

- **Simplicity** evaluates the ability of PRMs to detect redundancy in reasoning steps. Although redundant steps do not affect correctness, they increase computational costs and reduce efficiency. Additionally, simplifying the reasoning process enhances the clarity of the problem’s core and improves overall understandability.
- **Soundness** assesses the accuracy of the rewards produced by PRMs. As discussed in Section 1, errors in reasoning can vary in both causes and manifestations (Li et al., 2024b). Therefore, we evaluate not only the correctness of rewards but also the fine-grained performance across different error types and their nuances.
- **Sensitivity** measures PRMs’ robustness to details, such as critical conditions or implicit requirements. Sensitivity is vital for ensuring logical completeness and resilience to misleading information (Wen et al., 2025), contributing to the overall robustness of PRMs.

Each domain is further divided into detailed sub-categories for a more granular evaluation, which is discussed in detail below. The overall structure of PRMBENCH along with representative examples of each sub-category are illustrated in Figure 2, and the details of every evaluation category and sub-category are shown in Appendix A.1.

3.1.1 Simplicity

Specifically, the simplicity evaluation category is divided into two sub-categories: Non-Redundancy

	Overall	NR.	NCL.	ES.	SC.	DC.	CI.	PS.	DR.	MS.
Avg. Steps	13.4	15.3	10.3	13.8	14.2	13.3	14.2	12.7	13.4	14.1
Avg. Error Steps	2.1	2.0	2.8	2.8	1.6	1.8	1.7	2.5	2.3	0.0
Avg. First Error Step	7.8	7.8	4.9	8.0	9.1	6.8	11.4	6.2	8.3	N/A
Avg. Question Length	152.7	153.6	152.5	153.5	149.7	152.5	152.7	158.0	153.5	132.2
# of Instances	6216	758	758	757	758	757	757	756	750	165

Table 2: Statistics of PRMBENCH. NR., NCL., ES., SC., DC., CI., PS., DR., and MS. represent for Non-Redundancy, Non-Circular Logic, Empirical Soundness, Step Consistency, Domain Consistency, Confidence Invariance, Prerequisite Sensitivity, Deception Resistance, and Multi-Solution Consistency respectively.

and Non-Circular Logic, with detailed descriptions provided below:

Non-Redundancy evaluates the PRMs’ ability to identify redundancy within the reasoning process. Redundancy occurs when the reasoning includes unnecessary steps that do not contribute to the solution, making the process less concise and efficient. These steps can be removed without affecting the correctness of the final solution path.

Non-Circular Logic assesses the PRMs’ ability to detect circular reasoning within the process. Circular logic is a form of redundancy where the reasoning eventually loops back to a previous step, creating an infinite cycle. This sub-category is treated separately due to the frequent occurrence of circular logic in reasoning processes.

3.1.2 Soundness

We divide the Soundness category into four sub-categories due to its complexity: Empirical Soundness, Step Consistency, Domain Consistency, and Confidence Invariance. The definition of each sub-category is discussed below.

Empirical Soundness demands PRMs to detect the counterfactual mistakes within the reasoning process. A counterfactual step refers to a statement within a reasoning chain that contradicts established ground truth G .

Step Consistency expects PRMs to detect the step-wise contradiction, which means a conflict between a specific step and other steps within a reasoning path. Given a reasoning path $P = \{S_1, S_2, \dots, S_n\}$, a step contradiction exists if $S_i \perp S_j$, where $i, j \in [1, n]$ and $i \neq j$.

Domain Consistency requires PRMs to detect domain inconsistency mistakes, which is a special type of counterfactual. It refers to a step within the reasoning chain that uses a statement or theory

valid in other domains or cases but is not valid within the current reasoning chain.

Confidence Invariance demands PRMs to detect over-confident errors, a type of counterfactual where an incorrect statement is made with high confidence, contradicting established ground truth.

3.1.3 Sensitivity

This category includes three sub-categories: Prerequisite Sensitivity, Deception Resistance, and Multi-Solution Consistency, with detailed descriptions provided below.

Prerequisite Sensitivity requires PRMs to maintain sensitivity to missing conditions or prerequisite mistakes, which means a flaw in the reasoning chain where critical premises, assumptions, or necessary conditions are absent and this omission results in logical gaps, incomplete reasoning, or biased conclusions.

Deception Resistance demands PRMs to detect the deception or trap within a reasoning process, that is, statements that appear to be correct but are subtly altered to introduce inaccuracies while maintaining the illusion of correctness.

Multi-Solution Consistency expects PRMs to maintain consistency when faced with different solution paths of the same problem. Concretely, we utilize multiple correct reasoning processes of the same question to test whether the PRM can perform correctly.

3.2 Data Curation

We curate the dataset by extracting metadata and constructing test cases according to our category definitions. Detailed statistics of PRMBENCH are displayed in Table 2, with the curation procedure outlined below.

Meta Data Extraction Our metadata is built upon PRM800K (Lightman et al., 2023), which provides the questions (Q), ground truth answers (A), and ground truth step-level solution processes (S). We select completely correct solutions from both the training and test sets, filtering out low-quality instances to establish our ground truth answers.

Test Case Construction Each test case instance is represented as (Q', A, S') , where Q' denotes the test question and S' represents the test solution process, which may include errors. With class-specific prompts, as demonstrated in Appendix E.1, we query GPT-4o (OpenAI, 2024a) to modify the ground-truth reasoning process into versions containing erroneous steps. For the multi-solution, we leverage the newly proposed multi-step reasoning model QwQ¹ (Team, 2024a) to generate candidate answers for the given questions. These answers are then filtered to exclude unreasonable or incorrect ones, resulting in multi-solution reasoning processes for a single question.

3.3 Quality Control

To ensure a high-quality dataset, we implement a series of steps to filter out unqualified data and maintain data integrity. The specific procedures are outlined below:

Feature Filtering Our data curation procedure imposes strict structural requirements on the generated responses, where any outputs that do not satisfy these specifications cannot be considered valid for accurately assessing the performance of PRMs. However, even with detailed instructions, LLMs cannot consistently generate outputs that fully adhere to the required structure (Asai et al., 2024; Zeng et al., 2024a; Su et al., 2024a). To maintain high data quality, we define stringent filtering rules to exclude instances that fail to meet the necessary structural criteria. Detailed structural requirements are provided in Appendix E.1, and the full description of our data generation process can be found in the supplementary materials.

Human Verification Furthermore, to further ensure the quality of the data, we manually evaluate 10% of the total instances. We focus on two key qualities for each data instance: ❶ **Correctness of modification:** Whether the modifications made to the data instance are correct and reasonable. ❷

Difference in the modification: Whether the modified data instance differs from the original. We recruited five volunteers to evaluate our proposed PRMBENCH and observe over 92% qualification rate on the correctness metric and over 98% qualification rate on the difference metric. The details of human annotation are provided in Appendix A.3.1, and instructions for annotators are provided in Appendix D. This validation ensures the overall quality of our dataset and its suitability for studying process-level language reward models.

4 Experiments

4.1 Models

To provide a comprehensive evaluation of various models on PRMBENCH, we select a wide range of models, including open-source PRMs like Qwen-PRM (Zhang et al., 2025) and RLHFlowPRMs (Xiong et al., 2024), as well as LLMs prompted as critic models, such as o1-mini (OpenAI, 2024b) and DeepSeek R1 (Guo et al., 2025). A complete list of these models can be found in Appendix B.2. Additionally, we present the human evaluation results, with details available in Appendix A.3.2.

All PRMs and LLMs are evaluated on the complete PRMBENCH dataset, except for o1-mini and DeepSeek-R1, which are evaluated on a subset of PRMBENCH comprising 394 samples, proportionally selected to reflect the class distribution, in order to reduce evaluation costs.

Considering the complexity of the task, which involves question comprehension, evaluation of the provided processes, and adherence to format constraints, few-shot demonstration setups are employed to help the model adapt to the output format through In-Context Learning (ICL) examples. Specifically, we use two-shot examples when prompting general-purpose LLMs. The impact of few-shot settings is discussed in Section 5.4

4.2 Evaluation Metrics

Given our emphasis on evaluating the error detection capabilities, we use the negative F1 score as a metric for error detection performance. However, this metric may be affected by the inherent biases of models. To mitigate this and provide a unified, normalized score that reflects the overall competency of the evaluated model, following Zheng et al. (2024), we introduce a metric called PRMScore, defined formally in Equation 1.

¹Qwen/QwQ-32B-Preview: <https://huggingface.co/Qwen/QwQ-32B-Preview>

Model	Overall	Simplicity			Soundness					Sensitivity			
		NR.	NCL.	Avg.	ES	SC.	DC.	CI	Avg.	PS	DR.	MS.	Avg.
Human Performance	83.8	80.2	81.7	81.0	84.8	85.0	81.3	86.1	84.3	81.6	82.1	96.2	86.0
<i>Open-source Process Level Reward Models</i>													
Skywork-PRM-1.5B	61.1	52.0	56.4	54.2	64.8	64.9	63.3	66.5	64.9	57.5	63.3	91.1	70.7
Skywork-PRM-7B	65.1	<u>56.4</u>	62.8	59.6	69.4	67.1	<u>67.7</u>	69.9	68.5	60.9	65.8	93.2	73.3
Llemma-PRM800k-7B	52.0	49.3	53.4	51.4	56.4	47.1	46.7	53.3	50.9	51.0	53.5	93.6	66.0
Llemma-MetaMath-7B	50.5	50.2	50.5	50.3	51.9	47.6	44.4	52.1	49.0	50.5	51.3	96.0	66.0
Llemma-oprm-7B	50.3	48.7	49.3	49.0	54.2	46.8	44.5	53.5	49.8	49.2	51.3	91.8	64.1
MATHMinos-Mistral-7B	54.2	48.8	54.0	51.4	57.0	52.1	50.7	57.8	54.4	52.8	55.8	91.1	66.5
MathShepherd-Mistral-7B	47.0	44.0	50.3	47.1	49.4	44.5	41.3	47.7	45.7	47.2	48.6	86.1	60.7
ReasonEval-7B	60.1	61.0	50.1	<u>55.6</u>	62.1	65.9	61.5	66.0	63.9	55.7	58.0	99.5	71.1
ReasonEval-34B	60.5	54.8	48.1	51.5	66.4	60.3	57.8	67.5	63.0	57.7	64.3	97.2	73.1
RLHFlow-PRM-Mistral-8B	54.4	46.1	47.3	46.7	56.6	55.1	54.4	63.8	57.5	51.5	56.2	97.9	68.5
RLHFlow-PRM-Deepseek-8B	54.2	46.4	48.9	47.6	55.7	55.0	53.2	66.2	57.5	49.0	55.4	99.8	68.1
Qwen2.5-Math-PRM-7B	<u>65.5</u>	49.0	55.1	52.1	<u>71.8</u>	67.3	66.3	<u>78.5</u>	<u>71.0</u>	57.6	69.1	<u>99.7</u>	75.5
Qwen2.5-Math-PRM-72B	68.2	50.4	<u>58.8</u>	54.6	73.7	71.1	72.2	78.6	73.9	<u>60.3</u>	71.2	99.4	77.0
Pure-PRM-7B	65.3	49.2	55.2	52.2	71.1	<u>68.8</u>	64.0	76.9	70.2	60.3	<u>69.2</u>	98.0	<u>75.8</u>
Avg.	57.7	<u>50.5</u>	<u>52.9</u>	<u>51.7</u>	61.5	58.1	56.3	64.2	60.0	54.4	59.5	95.3	69.7
<i>Open LLMs, Prompted as Critic Models</i>													
MetaMath-7B	49.7	48.9	46.9	47.9	47.3	48.9	48.4	48.8	48.3	46.5	48.3	98.0	64.2
MetaMath-13B	49.4	50.3	44.4	47.3	47.8	47.4	49.4	48.1	48.2	49.0	48.1	99.5	65.5
Qwen2.5-Math-72B	57.4	55.3	54.9	55.1	55.5	<u>71.6</u>	58.1	59.1	61.1	47.4	53.8	100.0	67.1
QwQ-Preview-32B	<u>63.6</u>	<u>57.2</u>	<u>55.6</u>	<u>56.4</u>	<u>67.4</u>	72.3	<u>66.2</u>	<u>66.9</u>	<u>68.2</u>	<u>57.8</u>	<u>62.7</u>	100.0	<u>73.5</u>
R1-Distill-Llama3.1-70B	57.5	49.5	48.1	48.8	61.4	65.5	65.8	61.1	63.4	48.8	54.1	100.0	67.6
R1-Distill-Qwen-7B	52.6	32.9	37.9	35.4	47.3	54.1	48.4	48.0	49.4	45.6	46.8	100.0	64.1
DeepSeek-R1 [†]	67.8	63.0	62.7	62.9	68.2	68.5	73.5	75.4	71.4	63.3	68.0	100.0	77.1
Avg.	56.8	51.0	50.1	50.5	56.4	61.2	58.5	58.2	58.6	51.2	54.5	99.6	68.5
<i>Proprietary LLMs, Prompted as Critic Models</i>													
GPT-4o	66.8	57.0	62.4	59.7	72.0	<u>69.7</u>	70.7	71.1	70.9	62.5	65.7	99.2	75.8
o1-mini [†]	<u>68.8</u>	65.6	<u>63.7</u>	<u>64.6</u>	74.5	67.7	73.8	72.3	72.1	<u>61.8</u>	64.8	100.0	<u>75.5</u>
Gemini-2.0-flash-exp	66.0	<u>67.2</u>	58.1	62.7	70.4	65.7	66.0	67.3	67.3	61.8	66.2	98.2	75.4
Gemini-2.0-thinking-exp-1219	68.8	68.5	63.8	66.2	<u>72.9</u>	71.3	<u>71.0</u>	<u>71.8</u>	<u>71.8</u>	60.3	<u>65.7</u>	<u>99.8</u>	75.3
Avg.	67.6	64.6	62.0	63.3	72.4	68.6	70.4	70.7	70.5	61.6	65.6	99.3	75.5

Table 3: Performances comparison of popular models on PRMBENCH. The best performance for each category and task is in **bold**, while the second-best performance is underlined. [†]: To reduce costs, we evaluated only a subset of 394 samples for o1-mini and DeepSeek-R1.

$$PRM-Score = w_1 * F1_{neg} + w_2 * F1 \quad (1)$$

Where $F1$ and $F1_{neg}$ refer to $F1$ scores and negative $F1$ scores respectively. w_1 and w_2 are weights that are designed to maximize the differentiation between different models. The detailed evaluation procedure is provided in Appendix B.3. Besides, we also provide results of all evaluation categories in fine-grained metrics in Appendix B.4.

4.3 Main Results

The main results are shown in Table 3. Some observations can be summarized as follows:

The PRMs partially grasp multi-step process evaluation Our analysis indicates that, although Gemini-2-Thinking achieves the highest performance among all evaluated models, its score is still significantly lower than human performance (68.8

Model	Accuracy		PRM Score	Sim.
	Pos.	Neg.		
ReasonEval-7B	95.5	21.2	60.0	91.6
ReasonEval-34B	79.1	48.4	60.5	82.8
Skywork-7B	30.1	79.7	36.2	74.3
RLHFlow-DeepSeek-8B	95.0	13.0	54.2	95.0
GPT-4o	82.9	58.2	66.8	76.6
Gemini-2-thinking	89.0	49.8	68.8	82.0
Random	50.0	50.0	50.0	79.4

Table 4: Comparison of model performance on positive and negative test cases, along with their similarities.

vs. 83.8), highlighting substantial room for improvement in multi-step process evaluation. Some models even perform worse than random guessing, highlighting their limited reliability and potential training biases. Notably, the best open-source PRMs fail to match the performance of general-purpose proprietary LLMs, which suggests that even specifically trained PRMs still lag behind leading general-purpose models. We provide a detailed

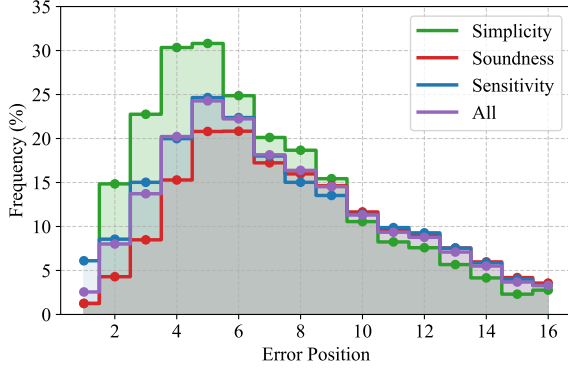


Figure 3: Distribution of error positions, truncated to 16 for better visualization, corresponding to the label field as shown in Figure 2.

error analysis in Appendix C.1.

Simplicity is more challenging for PRMs Our analysis highlights significant variations in model reasoning capabilities across evaluation categories. For instance, in the Sensitivity category, ReasonEval-34B performs relatively well, achieving an average score of 73.1. Especially in the Multi-Solutions sub-category, it excels with a PRM-Score of 97.2, approaching near-perfect classification accuracy. This suggests models perform relatively better on correct instance judgment. However, its performance declines markedly in more complex scenarios. In the **Simplicity** category, ReasonEval-34B’s PRMScore drops to 51.5, suggesting partially reliable performance.

Furthermore, to broaden the domain coverage of PRMBench, we additionally collect STEM-related data and construct PRMBENCH-STEM, which is designed to evaluate PRM performance across other domains. The data construction methodology and experimental results for PRMBench-STEM are provided in the appendix B.5.

5 Detailed Analysis

5.1 Inference Bias within PRMs

Takeaway 1. PRMs show a clear bias during evaluation, often favoring positive rewards.

As shown in Table 3, most open-source PRMs exhibit significant bias during evaluation, with some models performing worse than random guessing. suggesting the potential presence of bias within the inference procedure for our test cases. To validate this assumption, we compare the difference of models’ performance on positive and negative

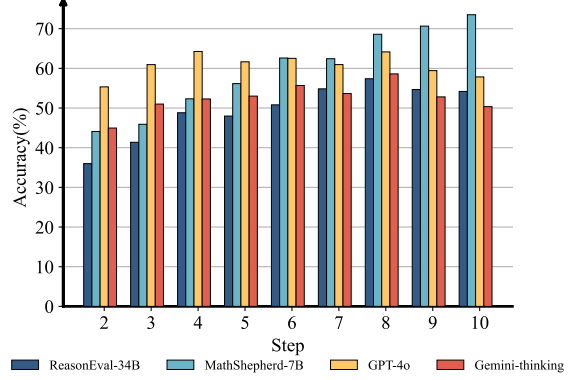


Figure 4: The models’ error-detection accuracy across different error steps, where step 1 and steps beyond 11 are truncated for improved visualization.

Model	0-shot	1-shot	2-shot
GPT-4o	68.1	68.2	66.8
Gemini-2-flash	65.3	64.9	66.0
Gemini-2-thinking	67.8	67.8	68.8

Table 5: The impact of ICL few-shot numbers on model performance. The number reported here is PRMScore.

instances. As shown in Table 4, **some models exhibit a clear bias during evaluation, often favoring positive rewards.** For instance, ReasonEval-7B and RLHFlow-DeepSeek-8B achieve over 95% accuracy on positive-labeled steps but only attain an average of 17% accuracy on negative-labeled steps. Although proprietary LLMs outperform open-source PRMs, they also exhibit bias with a comparatively milder reward tendency.

Additionally, to further investigate inference bias, we evaluate the reward similarity of models’ performance between completely correct reasoning processes and our test cases. The solution-level similarity is defined as $S = 100 - |Acc_{pos} - Acc_{neg}|$, where Acc denotes the average step accuracy within a solution. The results, shown in Table 4, reveal that certain models, such as ReasonEval-7B and RLHFlow-DeepSeek-8B, exhibit significantly higher similarity than the normal similarity score (79.4), showcasing potential limitations in differentiating positive and negative steps.

5.2 Performance across Different Steps

Takeaway 2. PRMs show a gradual improvement in performance as the position of the steps increases.

PRMBENCH includes a wide range of error step

Method	MATH	OlymBen	Avg.	PRMScore
Pass@8	96.2	79.8	88.0	-
Maj@8	71.8	40.3	56.1	-
SkyworkPRM-7B	90.0	60.1	75.1	65.1
LlemmaPRM-7B	87.4	58.3	72.8	52.0
MATHMinos-7B	88.3	59.1	73.7	54.2
MathShepherd-7B	88.6	60.0	74.3	47.0
ReasonEval-7B	87.0	58.4	72.7	60.1
RLHFlowPRM-8B	87.6	58.5	73.0	54.2
Qwen2.5-PRM-7B	88.0	58.7	73.4	65.5
Standard Dev (σ)	0.91	0.71	0.81	6.40
Somers' D	-0.05	0.05	-0.05	1.00

Table 6: Performance comparison on Best-of-8 using different PRMs. σ represents the standard deviation of model performances across all benchmarks. Somers' D refers to the Somers' D correlation between PRMScore and specific benchmarks.

positions. The distribution of error positions is illustrated in Figure 3. While differences exist across categories, the overall pattern remains consistent: all categories peak in frequency at step 5 and gradually decrease thereafter. This raises an interesting question: **Does the variation in step positions affect model performance?** To investigate, we focus on error steps to assess how erroneous step positions influence model accuracy. As depicted in Figure 4, proprietary LLMs maintain stable performance across different error step positions. In contrast, PRMs, including Math-Shepherd-7B and ReasonEval-7B, show a gradual improvement in performance as error step positions increase.

5.3 Impacts of ICL Settings

Takeaway 3. In-context learning has subtle impact on models' performance on PRMBENCH.

In this section, we investigate the impact of different ICL few-shot numbers on models' performance. We vary the number of ICL few-shots to 0, 1, and 2 to examine whether increasing the few-shot number enhances the performance of generative models prompted as critic models. As shown in Table 5, for the Gemini-series models, a subtle improvement in performance is observed with a few-shot setup. However, for GPT-4o, no significant improvement is detected, and in some cases, a larger few-shot number even results in a decline in performance. These findings suggest that a few-shot approach exerts only a subtle impact on model performance on PRMBENCH.

5.4 Comparison between BoN Evaluation and PRMBench

Takeaway 4. PRMs struggle with detecting false positives, exposing the potential for reward hacking.

We compare the results between our PRMBench and Best-of-N (BoN) evaluation to observe the correlation. Following Zhang et al. (2025); Yang et al. (2024a), we sampled eight responses (i.e., $N=8$) from Qwen-QwQ across multiple mathematical benchmarks, including GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021), Olympiad Bench (He et al., 2024a) and MMLU (Hendrycks et al., 2020). During evaluation, the PRMs are tasked with assigning a validity score to each step within every candidate response. The overall score for each candidate response is calculated by multiplying the individual step scores, as outlined in Lightman et al. (2023). We also provide majority voting as an baseline and pass@8 as the upper bound. The experiment setting and full BoN evaluation results are shown in Appendix C.2.

Although PRMs excel at selecting correct outcomes, they struggle with step-level reward hacking. As shown in Table 6, the average Somers' D correlation between PRMBench and BoN is only -0.05, highlighting the inconsistency between step-level and outcome-level evaluation. For instance, Math-Shepherd-7B achieves a PRMScore of 47.0 with 51.3% accuracy in false-positive scenarios within PRMBench, but outperforms most PRMs, including the state-of-the-art QwenPRM-7B, in the BoN evaluation (74.3 vs. 73.4). This inconsistency reveals that PRMs are suboptimal at detecting step-level errors and false positives, exposing potential reward hacking (Gao et al., 2025). Compared to BoN, PRMBENCH provides a better distinction between models, with a higher standard deviation (6.40 vs. 0.81), indicating its greater sensitivity to fine-grained differences in reasoning steps. Furthermore, based on the discoveries from PRMBench, we provide further discussions including several promising directions for future exploration, which are included in the appendix F.

6 Conclusion

In this paper, we investigate a crucial question: **Can existing PRMs detect various types of erroneous reasoning steps and provide reasonable rewards?** To address this, we introduce PRM-

BENCH, a benchmark characterized by its fine-grained evaluation categories and challenging requirements. We carefully curate 6,216 data samples with 83,456 step-level labels through LLMs and human filtering. PRMBENCH can be used to evaluate different process-labeling models, ensuring its general applicability. Through a comprehensive evaluation of existing PRMs and generative LLMs prompted as critic models, we can observe that PRMs exhibit partial capability in multi-step process evaluation, showcasing significant room for improvement. Furthermore, we highlight the critical need for detecting detailed error types and conducting comprehensive evaluations of PRMs. Despite these advances, enhancing the reward accuracy of PRMs and improving models’ reasoning abilities remain open research challenges. We encourage future work to leverage and expand upon PRMBENCH to address these issues.

7 Acknowledgement

We gratefully acknowledge the support and resources provided by the Shanghai Artificial Intelligence Laboratory, which played a vital role in enabling the successful completion of this research.

8 Limitations

There are still some limitations in our work, which are summarized below:

- While the PRMBench is large and comprehensive, comprising 6,216 samples and 83,456 step-level labels, a larger dataset could provide more robust evaluation and training opportunities. As the data construction method is flexible and data-agnostic, it can be adapted to more different data sources. We will continuously expand our dataset and explore PRM training in future versions.
- We evaluate the error-detection capabilities in terms of accuracy. However, a more detailed analysis, including the activation of the model’s neurons and hidden states (Zhang et al., 2023), would offer a deeper insight into improvements for PRMs. This limitation is not unique to our study and is common in most evaluations of Large Language Models.
- PRMBENCH currently focuses on textual reasoning. In future work, we will further explore its generalization to multimodal reasoning processes, such as those involving the integration of text and images (Su et al., 2025; Xia et al.,

2024a), to assess error detection capabilities in cross-modal scenarios.

References

- Akari Asai, Jacqueline He*, Rulin Shao*, Weijia Shi, Amanpreet Singh, Joseph Chee Chang, Kyle Lo, Luca Soldaini, Sergey Feldman, Tian, D’arcy Mike, David Wadden, Matt Latzke, Minyang, Pan Ji, Shengyan Liu, Hao Tong, Bohao Wu, Yanyu Xiong, Luke Zettlemoyer, Dan Weld, Graham Neubig, Doug Downey, Wen-tau Yih, Pang Wei Koh, and Hannaneh Hajishirzi. 2024. OpenScholar: Synthesizing scientific literature with retrieval-augmented language models. *Arxiv*.
- Jie Cheng, Lijun Li, Gang Xiong, Jing Shao, and Yisheng Lv. 2025. Pure: Prm is still effective and compute-efficient for llm math reasoning. <https://github.com/CJReinforce/PURE>.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- DeepMind. 2024. Gemini 2.0 flash experimental. <https://deepmind.google/technologies/gemini/flash/>. Accessed: 2024-12-25.
- Bofei Gao, Zefan Cai, Runxin Xu, Peiyi Wang, Ce Zheng, Runji Lin, Keming Lu, Junyang Lin, Chang Zhou, Wen Xiao, et al. 2024. Llm critics help catch bugs in mathematics: Towards a better mathematical verifier with natural language feedback. *CoRR*.
- Jiaxuan Gao, Shusheng Xu, Wenjie Ye, Weilin Liu, Chuyi He, Wei Fu, Zhiyu Mei, Guangju Wang, and Yi Wu. 2025. On designing effective RL reward at training time for LLM reasoning.
- Olga Golovneva, Moya Chen, Spencer Poff, Martin Corredor, Luke Zettlemoyer, Maryam Fazel-Zarandi, and Asli Celikyilmaz. 2022. Roscoe: A suite of metrics for scoring step-by-step reasoning. *arXiv preprint arXiv:2212.07919*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. 2024. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv preprint arXiv:2401.14196*.

- Yunzhuo Hao, Jiawei Gu, Huichen Will Wang, Linjie Li, Zhengyuan Yang, Lijuan Wang, and Yu Cheng. 2025. Can mllms reason in multimodality? emma: An enhanced multimodal reasoning benchmark. *arXiv preprint arXiv:2501.05444*.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. 2024a. Olympiad-bench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint arXiv:2402.14008*.
- Conghui He, Wei Li, Zhenjiang Jin, Chao Xu, Bin Wang, and Dahua Lin. 2024b. Opendatalab: Empowering general artificial intelligence with open datasets. *arXiv preprint arXiv:2407.13773*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- 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. *arXiv preprint arXiv:2103.03874*.
- Lei Li, Yuancheng Wei, Zhihui Xie, Xuqing Yang, Yifan Song, Peiyi Wang, Chenxin An, Tianyu Liu, Sujian Li, Bill Yuchen Lin, et al. 2024a. Vlerwardbench: A challenging benchmark for vision-language generative reward models. *arXiv preprint arXiv:2411.17451*.
- Xiaoyuan Li, Wenjie Wang, Moxin Li, Junrong Guo, Yang Zhang, and Fuli Feng. 2024b. Evaluating mathematical reasoning of large language models: A focus on error identification and correction. *arXiv preprint arXiv:2406.00755*.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*.
- Zicheng Lin, Zhibin Gou, Tian Liang, Ruilin Luo, Haowei Liu, and Yujiu Yang. 2024. Criticbench: Benchmarking llms for critique-correct reasoning. *arXiv preprint arXiv:2402.14809*.
- Yantao Liu, Zijun Yao, Rui Min, Yixin Cao, Lei Hou, and Juanzi Li. 2024. Rm-bench: Benchmarking reward models of language models with subtlety and style. *arXiv preprint arXiv:2410.16184*.
- Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. *arXiv preprint arXiv:2308.09583*.
- Nat McAleese, Rai Michael Pokorny, Juan Felipe Ceron Uribe, Evgenia Nitishinskaya, Maja Trebacz, and Jan Leike. 2024. Llm critics help catch llm bugs. *arXiv preprint arXiv:2407.00215*.
- Skywork o1 Team. 2024. [Skywork-o1 open series](https://huggingface.co/Skywork). <https://huggingface.co/Skywork>.
- OpenAI. 2024a. Gpt-4o system card. <https://cdn.openai.com/gpt-4o-system-card.pdf>. Accessed: 2024-09-26.
- OpenAI. 2024b. [Learning to reason with llms](https://openai.com/index/learning-to-reason-with-llms). <https://openai.com/index/learning-to-reason-with-llms/>.
- Yiwei Qin, Xuefeng Li, Haoyang Zou, Yixiu Liu, Shijie Xia, Zhen Huang, Yixin Ye, Weizhe Yuan, Hector Liu, Yuanzhi Li, et al. 2024. O1 replication journey: A strategic progress report—part 1. *arXiv preprint arXiv:2410.18982*.
- Xiaoye Qu, Yafu Li, Zhaochen Su, Weigao Sun, Jianhao Yan, Dongrui Liu, Ganqu Cui, Daizong Liu, Shuxian Liang, Junxian He, et al. 2025. A survey of efficient reasoning for large reasoning models: Language, multimodality, and beyond. *arXiv preprint arXiv:2503.21614*.
- Zhaochen Su, Linjie Li, Mingyang Song, Yunzhuo Hao, Zhengyuan Yang, Jun Zhang, Guanjie Chen, Jiawei Gu, Juntao Li, Xiaoye Qu, et al. 2025. Open-thinking: Learning to think with images via visual tool reinforcement learning. *arXiv preprint arXiv:2505.08617*.
- Zhaochen Su, Jun Zhang, Xiaoye Qu, Tong Zhu, Yanshu Li, Jiahuo Sun, Juntao Li, Min Zhang, and Yu Cheng. 2024a. Conflictbank: A benchmark for evaluating the influence of knowledge conflicts in llm. *arXiv preprint arXiv:2408.12076*.
- Zhaochen Su, Jun Zhang, Tong Zhu, Xiaoye Qu, Juntao Li, Min Zhang, and Yu Cheng. 2024b. Timo: Towards better temporal reasoning for language models. *arXiv preprint arXiv:2406.14192*.
- Qwen Team. 2024a. [Qwq: Reflect deeply on the boundaries of the unknown](#).
- ScalableMath Team. 2024b. [Easy-to-hard generalization models](https://huggingface.co/ScalableMath). <https://huggingface.co/ScalableMath>.
- Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. 2022. Solving math word problems with process-and outcome-based feedback. *arXiv preprint arXiv:2211.14275*.
- Peiyi Wang, Lei Li, Zhihong Shao, RX Xu, Damai Dai, Yifei Li, Deli Chen, Y Wu, and Zhifang Sui. 2023. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. *CoRR, abs/2312.08935*.

- Xueru Wen, Jie Lou, Yaojie Lu, Hongyu Lin, Xingyu, Xinyu Lu, Ben He, Xianpei Han, Debing Zhang, and Le Sun. 2025. [Rethinking reward model evaluation: Are we barking up the wrong tree?](#) In *The Thirteenth International Conference on Learning Representations*.
- Peng Xia, Kangyu Zhu, Haoran Li, Tianze Wang, Weijia Shi, Sheng Wang, Linjun Zhang, James Zou, and Huaxiu Yao. 2024a. Mmed-rag: Versatile multimodal rag system for medical vision language models. *arXiv preprint arXiv:2410.13085*.
- Shijie Xia, Xuefeng Li, Yixin Liu, Tongshuang Wu, and Pengfei Liu. 2024b. Evaluating mathematical reasoning beyond accuracy. *arXiv preprint arXiv:2404.05692*.
- Wei Xiong, Hanning Zhang, Nan Jiang, and Tong Zhang. 2024. An implementation of generative prm. <https://github.com/RLHFlow/RLHF-Reward-Modeling>.
- Yibo Yan, Shen Wang, Jiahao Huo, Hang Li, Boyan Li, Jiamin Su, Xiong Gao, Yi-Fan Zhang, Tianlong Xu, Zhendong Chu, et al. 2024. Errorradar: Benchmarking complex mathematical reasoning of multimodal large language models via error detection. *arXiv preprint arXiv:2410.04509*.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024a. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.
- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. 2024b. Qwen2. 5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhengguo Li, Adrian Weller, and Weiyang Liu. 2023. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*.
- Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng, Tanya Goyal, and Danqi Chen. 2024a. [Evaluating large language models at evaluating instruction following](#). In *The Twelfth International Conference on Learning Representations*.
- Zhongshen Zeng, Pengguang Chen, Shu Liu, Haiyun Jiang, and Jiaya Jia. 2023. Mr-gsm8k: A meta-reasoning benchmark for large language model evaluation. *arXiv preprint arXiv:2312.17080*.
- Zhongshen Zeng, Yinhong Liu, Yingjia Wan, Jingyao Li, Pengguang Chen, Jianbo Dai, Yuxuan Yao, Rongwu Xu, Zehan Qi, Wanru Zhao, Linling Shen, Jianqiao Lu, Haochen Tan, Yukang Chen, Hao Zhang, Zhan Shi, Bailin Wang, Zhijiang Guo, and Jiaya Jia. 2024b. [Mr-ben: A meta-reasoning benchmark for evaluating system-2 thinking in llms](#). *CoRR*, abs/2406.13975.
- Di Zhang, Jingdi Lei, Junxian Li, Xunzhi Wang, Yujie Liu, Zonglin Yang, Jiatong Li, Weida Wang, Suorong Yang, Jianbo Wu, et al. 2024a. Critic-v: Vlm critics help catch vlm errors in multimodal reasoning. *arXiv preprint arXiv:2411.18203*.
- Di Zhang, Jianbo Wu, Jingdi Lei, Tong Che, Jiatong Li, Tong Xie, Xiaoshui Huang, Shufei Zhang, Marco Pavone, Yuqiang Li, et al. 2024b. Llama-berry: Pairwise optimization for o1-like olympiad-level mathematical reasoning. *arXiv preprint arXiv:2410.02884*.
- 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. *arXiv preprint arXiv:2501.07301*.
- Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. 2023. Multimodal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*.
- Chujie Zheng, Zhenru Zhang, Beichen Zhang, Runji Lin, Keming Lu, Bowen Yu, Dayiheng Liu, Jingren Zhou, and Junyang Lin. 2024. Processbench: Identifying process errors in mathematical reasoning. *arXiv preprint arXiv:2412.06559*.
- Zihao Zhou, Shudong Liu, Maizhen Ning, Wei Liu, Jindong Wang, Derek F Wong, Xiaowei Huang, Qiufeng Wang, and Kaizhu Huang. 2024. Is your model really a good math reasoner? evaluating mathematical reasoning with checklist. *arXiv preprint arXiv:2407.08733*.

A Detailed Information for PRMBENCH

A.1 Evaluation Categories

In this section, We provide detailed information on our evaluation categories. The hierarchical categories, corresponding descriptions, and illustrations are shown in Figure 5. We carefully curated 6,216 data samples and 83,456 step-level labels. The benchmark spreads across three main evaluation categories: *simplicity*, *soundness*, and *sensitivity*. Among them, Simplicity comprises two sub-categories: non-redundancy and non-circular Logic. Soundness includes four main sub-categories: empirical soundness, step consistency, domain consistency, and confidence invariance. Finally, Sensitivity mainly evaluates models in three main parts: prerequisite sensitivity, deception resistance, and multi-solution consistency. The detailed descriptions and illustrations of each sub-category are shown in Figure 5.

A.1.1 Simplicity

Specifically, the Simplicity evaluation category is divided into two sub-categories: Non-Redundancy and Non-Circular Logic, with detailed descriptions provided below:

Non-Redundancy requires PRM to detect the redundancy within the reasoning procedure. The redundancy situation refers to a process that is not the most concise or efficient, as it includes one or more redundant steps that can be removed without affecting the correctness of the overall solution path. For example, as shown in Figure 5, if $A \rightarrow B$ represents a correct inference chain, the redundant reasoning procedure can be displayed as $A \rightarrow C \rightarrow B$, where C represents one or more redundant steps $C = \{c | c \text{ is redundant}\}$.

Non-Circular Logic In this sub-category, PRMs are required to detect the potential circular logic within the reasoning process. Circular logic is a specific form of redundancy, distinct from general redundancy, in that it finally loops back to a previous reasoning step. For example, as shown in Figure 2, if $A \rightarrow B$ represents a correct inference chain, circular logic can be formulated as $A \rightarrow C \rightarrow A \rightarrow B$, where the reasoning starts at step A , progresses through a sequence of steps, and ultimately loops back to A . We list Non-Circular Logic separately due to its common occurrence in reasoning processes.

A.1.2 Soundness

We divide the Soundness category into four sub-categories due to its complexity: Empirically Soundness, Step Consistency, Domain Consistency, and Confidence Invariance. The definition of each sub-category is discussed below.

Empirically Soundness demands PRM to detect the implicit counterfactual mistakes within the reasoning process. A counterfactual step refers to a statement within a reasoning chain that contradicts established ground truth G . Such contradictions can arise from relying on outdated theories, omitting critical constraints in theory, or incorporating erroneous assumptions.

Step Consistency expects PRM to detect the implicit step-wise contradiction, which means a conflict between a specific step and other steps within a reasoning path. Given a reasoning path $P = \{S_1, S_2, \dots, S_n\}$, a step contradiction exists if $S_i \perp S_j$, where $i, j \in [1, n]$ and $i \neq j$.

Domain Consistency Under this circumstance, PRMs are required to detect potential domain inconsistency mistakes, where domain inconsistency is a special type of counterfactual. It refers to a step within the reasoning chain that uses a statement or theory valid in other domains or cases but is not valid within the current reasoning chain.

Confidence Invariance demands the PRM to detect over-confident hallucinations, a type of counterfactual where an incorrect statement is made with unwarranted certainty, contradicting established ground truth.

A.1.3 Sensitivity

This category includes three sub-categories: Prerequisite Sensitivity, Deception Resistance, and Multi-Solution Consistency, with detailed descriptions provided below.

Prerequisite Sensitivity requires the PRM to maintain sensitivity to missing conditions or prerequisite mistakes, which means a flaw in the reasoning chain where critical premises, assumptions, or necessary conditions are absent. This omission results in logical gaps, incomplete reasoning, or biased conclusions. For example, when a missing condition occurs, the model is required to solve the problem through case analysis or further investigation. However, the answer becomes incorrect

Step Descriptions		GT	ReasonE	MathS	GPT-4o	Gemini
Question	Compute $(1 + \cos \frac{\pi}{8})(1 + \cos \frac{3\pi}{8})(1 + \cos \frac{5\pi}{8})(1 + \cos \frac{7\pi}{8})$.	-	-	-	-	-
Step 1	Let's call the expression we're trying to compute x .	✓	89.1	87.5	100.0	100.0
Step 2	Use the fact that $\cos(\pi - a) = -\cos a$.	✓	77.5	82.8	100.0	100.0
Step 3	Then we have $x = (1 + \cos \frac{\pi}{8})(1 + \cos \frac{3\pi}{8})(1 + \cos \frac{5\pi}{8})(1 + \cos \frac{7\pi}{8})$.	✓	71.4	76.7	100.0	100.0
Step 4	This equals $(1 + \cos \frac{\pi}{8})(1 + \cos \frac{3\pi}{8})(1 - \cos \frac{\pi}{8})(1 - \cos \frac{3\pi}{8})$.	✓	20.5	76.9	100.0	100.0
Step 5	By the difference of squares, this equals $(1 - \cos^2 \frac{\pi}{8})(1 - \cos^2 \frac{3\pi}{8})$.	✓	56.0	83.7	100.0	100.0
Step 6	And by the Pythagorean identity, this equals $(\sin^2 \frac{\pi}{8})(\sin^2 \frac{3\pi}{8})$.	✓	63.5	80.4	100.0	100.0
Step 7	Ok, now we just need to compute the sines.	✓	84.1	79.1	50.0	50.0
Step 8	We can use the formula for the sum of sines.	✗	73.6	74.1	0.0	0.0
Step 9	That formula is $\sin(a) + \sin(b) = 2 \sin(\frac{a+b}{2}) \cos(\frac{a-b}{2})$.	✗	79.1	78.5	0.0	100.0
Step 10	That means that $\sin \frac{\pi}{8} + \sin \frac{3\pi}{8} = 2 \sin(\frac{\pi/8+3\pi/8}{2}) \cos(\frac{\pi/8-3\pi/8}{2})$.	✗	67.4	78.0	0.0	0.0
Step 11	This equals $2 \sin \frac{\pi/2}{2} \cos(-\frac{\pi/4}{2})$.	✗	68.1	72.1	0.0	100.0
Step 12	This equals $2 \cdot 1 \cdot \frac{\sqrt{2}}{2}$.	✗	49.2	73.8	0.0	100.0
Step 13	This equals $\sqrt{2}$.	✗	35.1	73.8	0.0	100.0
Step 14	But this step doesn't help us compute the original product of sines.	✗	72.5	64.3	-50.0	100.0
Step 15	So, $(\sin^2 \frac{\pi}{8})(\sin^2 \frac{3\pi}{8})$ remains the same and $x = \frac{1}{8}$.	✗	6.3	35.8	0.0	-100.0
Step 16	Therefore, the mistake didn't change the value of x .	✓	22.6	43.5	-100.0	-100.0
Final Acc.	-	100	56.2	50.0	93.8	62.5
Reason	A counterfactual step was introduced in steps 8 through 13 by mistakenly using the formula for the sum of sines instead of the product of sines. This leads to incorrect intermediate calculations. However, due to fortunate errors, the end result ironically matches the correct answer in step 15.					

Table 7: An example of a data instance and error cases from PRMBENCH. The numbers reported are step-level validity scores generated by models. Scores and labels in **red** indicate negative samples, while those in **green** indicate positive samples. “GT” represents ground truth, while “ReasonE,” “MathS,” and “Gemini” correspond to ReasonEval-7B, Math-Shepherd-7B, and Gemini-2.0-flash-thinking-exp, respectively.

Abbr.	Full Name	Evaluation Category
NR.	Non-Redundancy	Simplicity
NCL.	Non-Circular Logic	Simplicity
ES.	Empirical Soundness	Soundness
SC.	Step Consistency	Soundness
DC.	Domain Consistency	Soundness
CI.	Confidence Invariance	Soundness
PS.	Prerequisite Sensitivity	Sensitivity
DR.	Deception Resistance	Sensitivity
MS.	Multi-Solution Consistency	Sensitivity

Table 8: The impact of ICL few-shot numbers on models’ final performance. The number reported here is PRMScore.

if the model overlooks the missing condition and proceeds with standard reasoning methods.

Deception Resistancy demands the PRM to detect the implicit deception or trap within a reasoning process, that is, statements that appear to be correct or align with ground truth but are subtly altered to introduce inaccuracies while maintaining the illusion of correctness.

Multi-Solution Consistency expects the PRM to maintain consistency when faced with different solution paths of the same problem. Concretely, to evaluate the sensitivity and the generalizability of PRMs, we utilize multiple correct reasoning processes of the same question to test whether the

PRM can perform correctly.

A.2 Examples For Different Evaluation Categories

In this section, we provide detailed examples of the various evaluation categories and their corresponding sub-categories. The data instance examples are displayed in Figure 7-18. All datasets used in this work are publicly available and have been released by their original creators, who are responsible for ensuring privacy protection. These datasets are utilized in accordance with their respective licenses and intended purposes, without introducing any harmful or sensitive information.

A.3 Human Annotation Settings

A.3.1 Quality Control

In the second stage of the quality control process, we recruited five volunteers, each holding a bachelor’s degree or an equivalent qualification, to assess the correctness and validity of PRMBENCH. To facilitate high-quality labeling, we utilize Label-LLM (He et al., 2024b) to help the data annotation procedure, as shown in Figure 6. The annotators’ instructions are shown in Appendix D.

A.3.2 Human Performance Evaluation

For human performance evaluation, we recruited three volunteers, each holding a bachelor’s degree or an equivalent qualification, to assist with data annotation. Following [Hao et al. \(2025\)](#), we randomly selected 50 instances from each sub-category, resulting in a mini-test set of 450 samples. Each annotator was responsible for three subsets, and the results are presented in Table 3.

B Detailed Experiment Results

B.1 Abbreviation Of Sub-Categories

The full names of abbreviations used in our tables are shown in Table 8.

B.2 Models

To provide a comprehensive evaluation of various models on PRMBENCH, we select a diverse set of models, including both open-source PRMs and different types of LLMs configured as critic models. Specifically, the open-source PRMs include Skywork-PRM-1.5B/7B ([o1 Team, 2024](#)), LlemmaPRMs ([Team, 2024b](#)), MathMinosPRM ([Gao et al., 2024](#)), MathShepherd-Mistral-7B ([Wang et al., 2023](#)), ReasonEval-7B/34B ([Xia et al., 2024b](#)), Pure-PRM-7B ([Cheng et al., 2025](#)) and Qwen-PRM-7B/72B ([Zhang et al., 2025](#)). Additionally, we evaluate state-of-the-art general-purpose LLMs, including the open-source MetaMath-7B/34B ([Yu et al., 2023](#)), Qwen2.5-Math-72B ([Yang et al., 2024b](#)), DeepSeek-R1 and its distill series ([Guo et al., 2025](#)), as well as closed-source LLMs such as GPT-4o ([OpenAI, 2024a](#)), Gemini-2.0-flash ([DeepMind, 2024](#)), and multi-step reasoning-enhanced LLMs like the o1 series models ([OpenAI, 2024b](#)) and Gemini-2-Thinking ([DeepMind, 2024](#)).

All models used in this work are publicly available and have been released by their original authors, who are responsible for ensuring privacy protection. These models are utilized in accordance with their respective licenses and intended purposes, without any modifications.

B.3 Evaluation Procedure

For each annotated question-solution pair, the reward models are tasked with evaluating the correctness and redundancy of each step, assigning a step-level validity score and a step-level redundancy score to each step. We subsequently utilize the specified threshold of each model to obtain the

Domain	# of Instances	Avg. Step Num	Avg. Error Num
Physics	1619	7.32	2.74
Chemistry	1543	7.74	2.78
Biology	2342	6.76	2.49

Table 9: Statistics of PRMBENCH-STEM.

prediction indicating whether the step is correct or redundant. This task is therefore framed as a binary classification problem. Thus we can utilize the evaluation metric defined in Section 4.2 to evaluate the performance of models on PRMBENCH.

B.4 Detailed Results of PRMBENCH

In addition to PRMScore displayed in Table 3, we also list the full results with different metrics across different sub-categories here. The detailed evaluation results are shown in Table 15-23.

B.5 Detailed Results of PRMBENCH-STEM

We extend our benchmark by collecting additional data from various scientific domains, including physics, chemistry, and biology, and construct PRMBENCH-STEM using a similar data curation methodology as described in the main paper. The statistics of PRMBENCH-STEM are presented in Table 9. We then evaluate several representative PRMs on PRMBENCH-STEM, and the results are summarized in Tables 10, 11, and 12 for the Biology, Chemistry, and Physics subsets, respectively. As shown, PRMs exhibit weaker performance on PRMBENCH-STEM compared to their performance on the original PRMBench. Notably, the Simplicity category remains the most challenging, consistent with our earlier observations.

C Details for Further Analysis

C.1 Error Analysis

A representative test case and the corresponding model performances are presented in Table 7. This example involves a counterfactual reasoning process, where steps eight through thirteen contain information that contradicts the correct computational principles and should be classified as “negative”. However, most models fail to identify these erroneous reasoning steps and assign relatively positive rewards, except for GPT-4o. While GPT-4o provides a relatively accurate reward, its judgments for key steps are only marginally negative, reflecting low confidence. This highlights a significant room for improvement in PRMs’ detailed error-detection capabilities.

Model	Overall	Simplicity			Soundness					Sensitivity		
		NR.	NCL.	Avg.	ES	SC.	DC.	CI	Avg.	PS	DR.	Avg.
Open-source Process Level Reward Models												
Skywork-PRM-1.5B	47.3	46.7	41.6	44.2	45.0	54.1	55.3	49.6	51.0	43.2	41.7	42.4
Skywork-PRM-7B	67.4	54.6	53.6	54.1	<u>71.6</u>	<u>76.2</u>	77.9	74.4	<u>75.0</u>	<u>64.0</u>	65.3	64.6
Llemma-PRM800k-7B	52.7	49.1	52.4	50.7	52.3	49.4	53.9	55.1	52.7	51.1	52.6	51.9
Llemma-MetaMath-7B	41.5	43.7	39.1	41.4	39.3	44.8	44.3	42.0	42.6	38.5	38.2	38.4
Llemma-oprm-7B	50.7	44.4	37.3	40.9	54.5	50.3	53.0	56.8	53.6	51.0	53.6	52.3
MATHMinos-Mistral-7B	47.1	43.0	37.9	40.4	48.0	51.6	47.1	53.6	50.1	45.1	45.3	45.2
MathShepherd-Mistral-7B	53.1	47.2	43.1	45.1	52.3	62.1	63.5	57.5	58.9	50.0	46.4	48.2
ReasonEval-7B	<u>67.5</u>	74.5	62.9	68.7	63.5	74.5	71.8	71.8	70.4	59.9	59.3	59.6
ReasonEval-34B	60.7	<u>68.2</u>	<u>61.9</u>	<u>65.0</u>	61.0	53.8	52.0	59.9	56.7	57.8	62.6	60.2
RLHFlow-PRM-Mistral-8B	45.5	43.0	37.4	40.2	44.5	53.1	49.0	50.3	49.2	42.6	42.2	42.4
RLHFlow-PRM-Deepseek-8B	47.0	42.9	37.4	40.1	44.7	59.0	57.7	52.6	53.5	41.6	41.6	41.6
Qwen2.5-Math-PRM-7B	65.8	42.9	49.0	46.0	69.6	75.1	72.5	<u>75.7</u>	73.2	63.9	<u>65.4</u>	<u>64.7</u>
Qwen2.5-Math-PRM-72B	73.7	42.8	50.6	46.7	82.3	79.6	<u>77.6</u>	82.9	80.6	75.5	76.6	76.1
Avg.	55.4	49.5	46.5	48.0	56.1	60.3	59.7	60.2	59.0	52.6	53.1	52.9
Open LLMs, Prompted as Critic Models												
MetaMath-7B	45.1	45.3	35.8	40.6	45.8	48.8	48.8	47.1	47.6	43.6	43.7	43.6
MetaMath-13B	43.6	44.7	35.8	40.2	42.2	46.3	47.1	45.0	45.2	42.4	41.1	41.7
Qwen2.5-Math-72B	42.3	49.5	35.6	42.5	38.0	56.1	45.7	42.8	45.6	37.0	35.6	36.3
QwQ-Preview-32B	54.5	54.8	43.3	49.0	53.6	63.9	<u>64.4</u>	58.5	<u>60.1</u>	51.8	48.3	50.1
R1-Distill-Llama3.1-70B	<u>53.8</u>	<u>49.8</u>	41.7	<u>45.7</u>	<u>53.4</u>	<u>63.6</u>	71.7	<u>57.8</u>	61.6	<u>46.7</u>	<u>44.5</u>	<u>45.6</u>
R1-Distill-Qwen-7B	50.4	39.9	<u>41.8</u>	40.9	46.4	48.7	50.9	47.4	48.3	45.6	43.3	44.5
Avg.	48.3	47.3	39.0	43.2	46.6	54.6	54.8	49.8	51.4	44.5	42.7	43.6

Table 10: Performance comparison of popular models on the **Biology** subset of PRMBENCH-STEM. The best performance for each category and task is in **bold**, while the second-best performance is underlined.

C.2 Details for BoN Evaluation

Following Zhang et al. (2025), we sample eight responses (i.e., $N=8$) from Qwen-QwQ-Preview 32B (Team, 2024a). During evaluation, the PRMs are tasked with assigning a validity score to each step within every candidate response. The overall score for each candidate response is calculated by multiplying the individual step scores, as outlined in Lightman et al. (2023). We then select the highest-ranked candidate response, compare it with the correct answer, and calculate the accuracy, which we refer to as $\text{prm}@8$. Additionally, we report the result of majority voting among the eight sampled responses ($\text{maj}@8$) as the baseline, and we define $\text{pass}@8$ as the proportion of test samples where any of the eight samplings lead to the correct final answer, which serves as the upper bound.

We conduct BoN evaluation across all models on GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021), Olympiad Bench (He et al., 2024a) and MMLU (Hendrycks et al., 2020). Due to the limit of the space, We omit some results in Table 6.

Thus we provide the full results of all PRMs in Table 13. Moreover, due to cost constraints, we select a subset of 200 samples for each benchmark and obtain a total of 800 samples to evaluate generative LLMs of their BoN performance. The results are shown in Table 14.

D Instructions for Human Annotators

D.1 Backgrounds

With the emergence of multi-step reasoning enhanced language models such as OpenAI o1, and Deepmind Gemini-thinking, these models demonstrate the ability to decompose complex problems and solve them step by step. However, while their solutions often appear correct, they may contain errors in understanding, calculation, or reasoning logic, which is also known as false positive situations. A popular way to evaluate the results generated by these models is by utilizing process-level reward models (PRMs). Nevertheless, PRMs are fallible and not always correct. Existing benchmarks are not adequate for evaluating PRMs on

Model	Overall	Simplicity			Soundness					Sensitivity		
		NR.	NCL.	Avg.	ES	SC.	DC.	CI	Avg.	PS	DR.	Avg.
Open-source Process Level Reward Models												
Skywork-PRM-1.5B	48.1	45.9	43.9	44.9	46.3	52.5	54.0	51.3	51.0	45.3	43.3	44.3
Skywork-PRM-7B	64.3	51.2	51.1	51.1	65.2	<u>70.4</u>	<u>71.6</u>	73.3	70.1	64.2	64.6	64.4
Llemma-PRM800k-7B	52.6	47.1	53.2	50.2	56.5	48.8	47.9	55.7	52.2	52.6	55.2	53.9
Llemma-MetaMath-7B	45.2	45.3	41.5	43.4	44.9	48.2	46.3	47.6	46.7	42.3	42.8	42.5
Llemma-oprm-7B	48.3	44.7	38.1	41.4	53.5	46.6	41.6	56.2	49.5	49.0	52.4	50.7
MATHMinos-Mistral-7B	49.6	44.0	38.8	41.4	50.8	51.9	46.3	58.2	51.8	49.7	50.4	50.0
MathShepherd-Mistral-7B	52.5	47.8	43.4	45.6	55.7	56.3	54.3	58.4	56.2	50.2	49.5	49.8
ReasonEval-7B	62.9	69.4	60.1	64.7	59.2	66.3	64.1	70.3	65.0	57.4	55.5	56.4
ReasonEval-34B	63.2	<u>63.5</u>	<u>55.8</u>	<u>59.6</u>	65.6	58.6	57.8	66.9	62.2	61.9	65.6	63.7
RLHFlow-PRM-Mistral-8B	45.6	43.3	38.5	40.9	44.2	50.9	47.7	51.2	48.5	44.1	42.3	43.2
RLHFlow-PRM-Deepseek-8B	47.3	43.5	38.3	40.9	46.9	51.8	50.5	56.3	51.4	45.4	43.1	44.2
Qwen2.5-Math-PRM-7B	<u>66.3</u>	43.5	49.9	46.7	<u>70.9</u>	66.8	68.8	<u>80.0</u>	<u>71.6</u>	<u>67.0</u>	<u>68.1</u>	<u>67.5</u>
Qwen2.5-Math-PRM-72B	71.4	44.8	53.2	49.0	77.3	71.6	74.6	81.0	76.1	73.1	77.5	75.3
Avg.	55.2	48.8	46.6	47.7	56.7	57.0	55.8	62.0	57.9	54.0	54.6	54.3
Open LLMs, Prompted as Critic Models												
MetaMath-7B	45.6	44.2	37.5	40.9	46.0	48.5	48.0	47.3	47.4	45.6	<u>45.5</u>	<u>45.5</u>
MetaMath-13B	44.3	42.9	37.0	40.0	44.7	48.6	47.0	45.3	46.4	45.2	42.2	43.7
Qwen2.5-Math-72B	43.2	50.5	34.8	<u>42.6</u>	37.6	<u>61.7</u>	43.8	43.3	46.6	37.5	36.6	37.0
R1-Distill-Llama3.1-70B	<u>50.6</u>	<u>45.9</u>	41.0	43.5	48.2	65.7	58.1	57.8	57.5	<u>46.7</u>	42.8	44.8
R1-Distill-Qwen-7B	51.8	35.6	<u>41.0</u>	38.3	<u>47.5</u>	51.0	<u>50.1</u>	<u>49.2</u>	<u>49.5</u>	47.1	46.2	46.7
Avg.	47.1	43.8	38.2	41.0	44.8	55.1	49.4	48.6	49.5	44.4	42.7	43.5

Table 11: Performance comparison of popular models on the **Chemistry** subset of PRMBENCH-STEM. The best performance for each category and task is in **bold**, while the second-best performance is underlined.

different error types. Therefore, we are building a comprehensive evaluation benchmark for PRMs that can have a fine-grained detection of PRMs.

D.2 Task Definition

We begin by collecting completely correct multi-solution data and leveraging state-of-the-art LLMs to introduce various types of errors into these correct solutions, thereby generating our test cases. The detailed error types are described in Section 3. All synthesized data instances undergo an initial filtering process based on specific features.

Your task is to identify whether the modification taken is reasonable and whether the modified data instance is different from the original data instance.

Sub-task 1

The first sub-task is a binary classification task whose options include yes and no. Your task is to decide whether the modified step-by-step solution generated by LLMs is reasonable. The word “reasonable” has two aspects for evaluation.

- The modified process generated by LLMs seems like a possible solution path that could happen.

- The modified process generated by LLMs is exactly wrong and the type of error is suitable for the current “classification”.

Please assign a “yes” for this sub-task if both of the answers to the above two questions are “yes”. Otherwise, assign a “no” for this sub-task.

D.2.1 Sub-task 2

The second sub-task is a binary classification task whose options include yes and no. Your task is to decide whether the modified step-by-step solution generated by LLMs is different from the original solution process. The word “different” means the modified solution process is logically different from the original one, or there exist different statements compared to the original process.

Please assign a “yes” to this sub-task if your answer to the above question is “yes”. Otherwise, assign a “no” for this sub-task.

D.2.2 Error Types

Redundancy refers to a process that is not the most concise or efficient, as it includes one or more

Model	Overall	Simplicity			Soundness					Sensitivity		
		NR.	NCL.	Avg.	ES	SC.	DC.	CI	Avg.	PS	DR.	Avg.
Open-source Process Level Reward Models												
Skywork-PRM-1.5B	50.3	44.7	41.8	43.3	51.8	57.7	55.2	56.1	55.2	46.8	46.8	46.8
Skywork-PRM-7B	65.2	51.1	51.9	51.5	68.3	<u>70.3</u>	72.7	76.4	71.9	64.8	65.1	65.0
Llemma-PRM800k-7B	53.8	49.1	52.0	50.6	57.3	49.6	49.4	57.3	53.4	53.8	58.6	56.2
Llemma-MetaMath-7B	46.9	46.4	42.0	44.2	46.5	47.2	44.2	51.7	47.4	47.8	45.4	46.6
Llemma-oprm-7B	47.1	42.8	39.7	41.2	52.2	42.7	40.9	54.3	47.5	48.1	51.6	49.8
MATHMinos-Mistral-7B	46.5	43.6	38.9	41.2	46.9	47.8	44.2	57.8	49.2	45.3	43.3	44.3
MathShepherd-Mistral-7B	53.5	48.0	44.8	46.4	55.9	56.3	54.8	62.8	57.4	48.6	53.5	51.1
ReasonEval-7B	62.5	71.5	57.2	64.4	60.1	66.9	62.2	70.6	65.0	56.9	53.9	55.4
ReasonEval-34B	64.2	<u>68.5</u>	52.0	<u>60.2</u>	68.5	59.3	56.4	69.8	63.5	63.7	66.6	65.2
RLHFlow-PRM-Mistral-8B	46.5	42.7	37.5	40.1	46.0	51.0	50.5	55.7	50.8	44.2	42.2	43.2
Qwen2.5-Math-PRM-7B	<u>66.9</u>	43.7	49.3	46.5	<u>73.6</u>	66.4	70.3	<u>80.0</u>	<u>72.6</u>	<u>69.6</u>	<u>68.8</u>	<u>69.2</u>
Qwen2.5-Math-PRM-72B	71.1	44.2	<u>56.2</u>	50.2	78.9	71.1	<u>71.1</u>	80.3	75.4	72.9	77.6	75.2
Avg.	55.6	45.9	46.3	46.1	58.1	56.8	55.7	64.2	58.7	54.5	55.2	54.8
Open LLMs, Prompted as Critic Models												
MetaMath-7B	45.6	43.5	36.9	40.2	<u>47.9</u>	48.9	<u>49.1</u>	48.1	48.5	<u>47.3</u>	<u>43.6</u>	<u>45.4</u>
MetaMath-13B	43.8	42.7	36.6	39.7	44.5	46.7	48.0	46.3	46.4	41.5	41.9	41.7
Qwen2.5-Math-72B	44.8	50.6	<u>37.8</u>	44.2	39.8	64.5	45.5	47.5	<u>49.3</u>	37.9	36.2	37.1
R1-Distill-Llama3.1-70B	53.3	<u>44.5</u>	40.4	<u>42.5</u>	53.0	<u>63.8</u>	59.4	64.6	60.2	49.6	48.3	48.9
R1-Distill-Qwen-7B	<u>51.0</u>	29.7	34.8	32.3	44.1	48.1	45.7	<u>48.9</u>	46.7	44.0	40.7	42.3
Avg.	47.7	42.2	37.3	39.8	45.9	54.4	49.5	51.1	50.2	44.0	42.2	43.1

Table 12: Performance comparison of popular models on the **Physics** subset of PRMBENCH-STEM. The best performance for each category and task is in **bold**, while the second-best performance is underlined.

Method	GSM8K	MATH	OlymBen	MMLU	Avg.	PRMScore
Pass@8	98.9	96.2	79.8	96.7	92.9	-
Maj@8	95.4	71.8	40.3	85.8	73.3	-
Skywork-PRM-1.5B	96.7	89.2	58.8	90.2	83.7	61.1
Skywork-PRM-7B	97.1	90.0	60.1	90.3	84.4	65.1
Llemma-PRM800k-7B	96.0	87.4	58.3	90.0	82.9	52.0
Llemma-MetaMath-7B	96.0	88.2	58.6	90.0	83.2	50.5
Llemma-oprm-7B	96.4	86.6	58.0	89.9	82.7	50.3
MATHMinos-Mistral-7B	95.8	88.3	59.1	89.1	83.1	54.2
MathShepherd-Mistral-7B	96.6	88.6	60.0	90.0	83.8	47.0
ReasonEval-7B	96.4	87.0	58.4	90.0	82.9	60.1
RLHFlow-PRM-Mistral-8B	96.5	88.4	59.1	90.4	83.6	54.4
RLHFlow-PRM-Deepseek-8B	96.4	87.6	58.5	90.2	83.2	54.2
ReasonEval-34B	96.4	86.9	56.9	90.0	82.6	60.5
Qwen2.5-Math-PRM-7B	96.7	88.0	58.7	89.8	83.3	65.5
Qwen2.5-Math-PRM-72B	96.7	89.3	60.5	90.0	84.2	68.2
Standard Deviation (σ)	0.35	0.97	0.93	0.29	0.53	6.41
Somers' D correlation	0.50	0.26	0.22	0.19	0.28	1.00

Method	GSM8K	MATH	OlymBen	MMLU	Avg.	PRMScore
Pass@8	99.0	96.0	77.0	94.0	91.5	-
Maj@8	96.5	68.0	41.0	86.0	72.9	-
QwQ-Preview-32B	96.5	83.5	56.5	87.5	81.0	63.6
MetaMath-7B	95.5	82.0	58.0	88.0	80.9	49.7
MetaMath-13B	95.0	84.5	57.5	86.5	80.9	49.4
MetaMath-70B	96.5	85.0	58.5	86.0	81.5	45.9
Qwen2.5-Math-7B	96.0	85.5	59.0	86.5	81.8	49.2
Qwen2.5-Math-72B	97.5	82.0	53.5	88.5	80.4	57.4
R1-Distill-Llama3.1-8B	96.5	84.5	53.5	87.5	80.5	52.7
R1-Distill-Llama3.1-70B	95.0	83.0	59.5	86.5	81.0	57.5
R1-Distill-Qwen-7B	97.0	86.5	59.0	86.5	82.2	52.6
R1-Distill-Qwen-32B	97.0	82.0	60.5	84.5	81.0	60.2
WizardMath-7B	96.5	82.5	60.0	87.0	81.5	49.2
Gemini-2.0-flash-exp	97.0	81.5	56.5	86.5	80.4	66.0
Gemini-2.0-thinking-exp-1219	98.0	87.5	60.5	89.5	83.9	68.8
Standard Deviation (σ)	0.87	1.83	2.25	1.19	0.91	7.02
Somers' D correlation	0.36	-0.21	0.03	0.24	-0.15	1.00

Table 13: Performance comparison on Best-of-8 using different **PRMs**. σ represents the standard deviation of model performances across all benchmarks. Somers' D refers to the Somers' D correlation between PRMScore and specific benchmarks.

Table 14: Performance comparison on Best-of-8 using different **LLMs as a Judge**. σ represents the standard deviation of model performances across all benchmarks. Somers' D refers to the Somers' D correlation between PRMScore and specific benchmarks.

redundant steps that can be removed without affecting the correctness of the overall solution path. For example, if $A \rightarrow B$ represents a correct inference chain, your task is to introduce one or more redundant steps $C = \{c | c \text{ is redundant}\}$ and reformulate the solution chain as $A \rightarrow C \rightarrow B$.

Circular logic is a specific form of redundancy, characterized by a reasoning chain that starts at a step S , progresses through a sequence of steps, and ultimately loops back to S . Symbolically, this can be expressed as $S \rightarrow A \rightarrow B \rightarrow S$, where S , A , and B represent individual reasoning steps. Your task is to modify the reasoning process to introduce such circular logic.

Counterfactual A counterfactual step refers to a statement within a reasoning chain that contradicts ground truth or established theories. Such contradictions can arise from relying on outdated theories, omitting critical constraints in a theory, or incorporating erroneous assumptions. Your task is to modify the reasoning process to introduce such counterfactual steps.

Step contradiction refers to a conflict between a specific step and other steps within a reasoning path. Given a reasoning path $P = S_1, S_2, \dots, S_n$, a step contradiction exists if $S_i \perp S_j$, where $i, j \in [1, n]$ and $i \neq j$. Your task is to modify the reasoning process to introduce such step contradiction steps.

Domain inconsistency is a special type of counterfactual. It refers to a step or a few steps within the reasoning chain that uses a statement or theory valid in other domains or cases but is not valid within the current reasoning chain. Your task is to modify the reasoning process to introduce steps for such domain inconsistency.

Confident hallucination is a special type of counterfactual. It refers to a statement within the reasoning chain that contradicts established ground truth and is presented with an overly confident tone. In other words, it involves stating an incorrect statement with unwarranted certainty. Your task is to modify the reasoning process to introduce such confident hallucination steps.

Missing condition or prerequisite refers to a flaw in the reasoning chain where critical premises, assumptions, or necessary conditions are absent. This omission results in logical gaps, incomplete reasoning, or biased conclusions. For example,

when a missing condition occurs, the model must solve the problem through case analysis or further investigation. However, the answer becomes incorrect if the model overlooks the missing condition and proceeds with standard reasoning methods. Your task is to modify the reasoning process to introduce such missing condition errors.

Deception or traps refer to statements that appear to be correct or align with ground truth but are subtly altered to introduce inaccuracies while maintaining the illusion of correctness. Your task is to modify the reasoning process to introduce such deception or trap error steps.

E Prompts

E.1 Prompts For Generating Data

As introduced in Section 3.2, we query GPT-4o (OpenAI, 2024a) to synthesize the metadata at the very first step of our test case construction procedure. To better prompt LLMs to generate high-quality data instances, we carefully designed our prompts, which are displayed in Figure 19-22. We display only one example here due to limitations in space. And the prompts can be found in our supplementary materials.

E.2 Prompts For Evaluating Generative LLMs

As introduced in Section 4.1, we prompt some state-of-the-art generative LLMs as critic models to evaluate their rewarding capabilities on PRMBENCH. To make a fair comparison between different models, we carefully design the prompts and utilize a unified prompt to query them. The prompt used is displayed in Figure 23 and 24.

F Further Discussion

Inspired by the results and discoveries on PRMBENCH, we further propose several promising directions for future research, which we hope can offer valuable insights and contribute to the advancement of the research community.

Anti-redundancy training: As stated in Section 4.3, our work highlights a specific weakness of current PRMs in identifying redundant reasoning steps. To mitigate this, one possible approach is to modify the label distribution during training. PRM training data is typically labeled as Correct, Neutral, or Incorrect, where the Neutral label often corresponds to redundant steps. By reducing the proportion of

Neutral samples in the training data, we can train PRMs with stronger anti-redundancy capabilities.

Contrastive training: A high-quality data curation pipeline is introduced in Section 3.1, which can also be adapted to curate training samples labeled with fine-grained error types. By leveraging contrastive learning or preference alignment with the curated data, the error sensitivity and detection capabilities of PRMs can be further improved.

Step-level evaluation for LLMs: As introduced in Section 5.4, the inconsistency between PRM-Bench and BoN evaluation reveals the false positive situation and the risk of reward hacking within LM post-training. Therefore, traditional outcome-level label-based evaluation is not enough, highlighting the need for a comprehensive step-level evaluation of LLM’s reasoning procedure.

Model Name	PRMScore	F1	Negative F1	Acc	Positive Acc	Negative Acc	First	similarity
<i>Open-source Process Level Reward Models</i>								
Skywork-PRM-1.5B	61.1	88.6	33.6	80.6	88.8	33.2	44.7	91.7
Skywork-PRM-7B	65.1	89.2	<u>40.9</u>	81.7	88.5	42.7	56.6	90.1
Llemma-PRM800k-7B	52.0	75.7	28.3	63.7	66.4	<u>48.5</u>	22.2	79.5
Llemma-MetaMath-7B	50.5	80.4	20.7	68.5	75.6	27.7	15.1	84.2
Llemma-oprm-7B	50.3	77.3	23.3	64.9	69.9	36.1	16.6	83.7
MATHMinos-Mistral-7B	54.2	79.2	29.1	67.9	72.8	41.7	38.0	84.6
MathShepherd-Mistral-7B	47.0	64.9	29.2	53.0	51.5	61.1	<u>54.6</u>	83.5
ReasonEval-7B	60.1	90.8	29.3	83.8	95.5	21.2	30.3	91.6
ReasonEval-34B	60.5	83.8	37.2	74.2	79.1	48.4	50.8	82.8
RLHFlow-PRM-Mistral-8B	54.4	87.7	21.1	78.8	90.2	17.9	22.1	92.8
RLHFlow-PRM-Deepseek-8B	54.2	89.9	18.6	82.0	95.0	13.0	17.0	95.0
Qwen2.5-Math-PRM-7B	<u>65.5</u>	91.5	39.4	85.1	<u>95.4</u>	30.6	37.8	89.0
Qwen2.5-Math-PRM-72B	68.2	<u>91.4</u>	45.1	<u>85.1</u>	93.8	38.7	48.5	86.8
Pure-PRM-7B	65.3	90.1	40.5	83.0	91.8	36.6	41.8	86.6
Avg.	57.7	84.3	31.2	75.2	82.4	35.5	35.4	87.3
<i>Open LLMs, Prompted as Critic Models</i>								
MetaMath-7B	49.7	88.2	11.2	79.1	90.9	9.1	7.6	91.6
MetaMath-13B	49.4	89.6	9.1	81.4	94.4	6.3	4.8	94.8
Qwen2.5-Math-72B	57.4	<u>90.3</u>	24.4	<u>82.9</u>	<u>96.9</u>	15.9	19.8	92.4
QwQ-Preview-32B	<u>63.6</u>	87.6	<u>39.6</u>	79.4	89.2	<u>36.4</u>	<u>40.2</u>	83.1
R1-Distill-Llama3.1-70B	57.5	91.4	23.5	84.6	97.4	15.2	19.0	93.8
R1-Distill-Qwen-7B	52.6	83.2	22.0	72.3	80.5	26.0	23.9	79.3
DeepSeek-R1 [†]	67.8	87.2	48.4	79.5	83.1	60.5	63.7	77.9
Avg.	56.8	88.2	25.4	79.9	90.3	24.2	25.6	87.5
<i>Proprietary LLMs, Prompted as Critic Models</i>								
GPT-4o	66.8	86.9	46.7	79.0	82.9	58.2	64.4	76.6
o1-mini [†]	<u>68.8</u>	<u>89.2</u>	48.3	<u>82.1</u>	<u>86.9</u>	55.4	56.4	80.4
Gemini-2.0-flash-exp	66.0	86.5	45.5	78.4	82.3	<u>57.2</u>	<u>64.3</u>	80.0
Gemini-2.0-thinking-exp-1219	68.8	89.7	<u>47.8</u>	82.8	89.0	49.8	57.0	82.0
Avg.	67.6	88.1	47.1	80.6	85.3	55.2	60.5	79.7

Table 15: A performance comparison of popular models across detailed metrics in **ALL** categories of PRMBENCH. The best performance for each metric is highlighted in **bold**, while the second-best performance is underlined. [†]: To reduce costs, we evaluated only a subset of 394 samples for o1-mini and Deepseek-R1.

	Categories	Descriptions	Illustration
Simplicity	Non-Redundancy	Non-Redundancy requires the PRM to detect redundant situations, that is, a process that includes one or more redundant steps that can be removed without affecting the correctness of the overall solution.	
	Non-Circular Logic	Circular logic is a specific form of redundancy, characterized by a reasoning chain that starts at a step A, progresses through a sequence of steps, and ultimately loops back to A.	
Soundness	Empirical Soundness	Empirical Soundness requires the prm to detect a counterfactual step, concretely, a statement within a reasoning chain that contradicts established ground truth.	
	Step Consistency	Step Consistency requires the PRM to detect a conflict between a specific step and other steps within a reasoning path.	
	Domain Consistency	Domain Consistency requires the PRM to maintain robust when faced with a step that uses a statement or theory valid in other domains or cases but is not valid within the current reasoning scene.	
	Confidence Invariance	Confidence Invariance requires the PRM to maintain invariant when faced with a statement within the reasoning chain that contradicts established ground truth and is presented with an overly confident tone.	
Sensitivity	Prerequisite Sensitivity	Prerequisite Sensitivity requires the PRM to detect the flaw in the reasoning chain where critical premises, assumptions, or necessary conditions are absent and the absense would cause errors.	
	Deception Resistance	Deception Resistance requires the PRM to detect statements that appear to be correct or align with ground truth but are subtly altered to introduce inaccuracies while maintaining the illusion of correctness.	
	Multi-Solution Consistency	Multi-Solution Consistency requires the PRM to maintain consistent when faced with different solution paths of the same problem.	

Figure 5: Examples of an in-depth evaluation of PRMBench.

Model Name	PRMScore	F1	Negative F1	Acc	Positive Acc	Negative Acc	First	similarity
<i>Open-source Process Level Reward Models</i>								
Skywork-PRM-1.5B	52.0	88.4	15.6	79.6	88.8	15.1	18.7	95.0
Skywork-PRM-7B	<u>56.4</u>	89.0	<u>23.9</u>	80.7	88.8	24.2	31.3	94.5
Llemma-PRM800k-7B	49.3	77.7	20.9	65.2	69.2	<u>36.9</u>	20.5	86.4
Llemma-MetaMath-7B	50.2	85.3	15.0	75.0	83.2	17.6	15.3	94.4
Llemma-oprm-7B	48.7	80.3	17.2	68.1	74.1	26.5	16.0	92.0
MATHMinos-Mistral-7B	48.8	82.7	14.9	71.2	79.2	19.0	20.8	95.7
MathShepherd-Mistral-7B	44.0	67.7	20.2	54.0	55.6	43.7	39.7	93.8
ReasonEval-7B	61.0	91.5	30.5	84.9	94.0	25.1	<u>32.8</u>	89.8
ReasonEval-34B	54.8	86.7	22.9	77.4	85.3	25.3	31.6	83.8
RLHFlow-PRM-Mistral-8B	46.1	89.1	3.2	80.3	92.3	2.4	3.4	98.5
RLHFlow-PRM-Deepseek-8B	46.4	91.0	1.9	83.5	96.1	1.2	1.5	99.0
Qwen2.5-Math-PRM-7B	49.0	92.1	5.9	85.4	98.0	3.5	3.3	98.2
Qwen2.5-Math-PRM-72B	50.4	<u>91.9</u>	8.8	<u>85.2</u>	<u>97.4</u>	5.4	6.2	98.2
Pure-PRM-7B	49.2	90.8	7.6	83.3	95.2	5.2	5.4	96.2
Avg.	50.5	86.0	14.9	76.7	85.5	17.9	17.6	94.0
<i>Open LLMs, Prompted as Critic Models</i>								
MetaMath-7B	48.9	80.7	17.0	68.7	74.9	25.5	23.2	76.3
MetaMath-13B	50.3	86.5	14.0	76.7	85.7	14.9	11.8	85.4
Qwen2.5-Math-72B	55.3	<u>90.3</u>	20.3	<u>82.7</u>	<u>93.9</u>	15.4	17.0	90.1
QwQ-Preview-32B	<u>57.2</u>	85.8	<u>28.7</u>	76.3	84.3	31.4	33.9	81.9
R1-Distill-Llama3.1-70B	49.5	92.5	6.5	86.1	98.8	3.6	5.4	98.0
R1-Distill-Qwen-7B	32.9	43.9	21.9	34.7	29.4	70.1	68.6	31.5
DeepSeek-R1 [†]	63.0	86.9	39.2	78.4	82.2	<u>53.3</u>	<u>62.5</u>	76.5
Avg.	51.0	80.9	21.1	71.9	78.5	30.6	31.8	77.1
<i>Proprietary LLMs, Prompted as Critic Models</i>								
GPT-4o	57.0	77.8	36.3	67.0	66.5	70.4	77.1	68.9
o1-mini [†]	65.6	90.8	40.4	84.1	90.4	41.5	45.8	84.1
Gemini-2.0-flash-exp	<u>67.2</u>	91.5	<u>42.9</u>	85.1	91.8	41.7	49.7	82.4
Gemini-2.0-thinking-exp-1219	68.5	<u>91.4</u>	45.6	<u>85.1</u>	<u>90.9</u>	<u>47.1</u>	<u>56.5</u>	84.7
Avg.	64.6	87.9	41.3	80.4	84.9	50.2	57.3	80.0

Table 16: A performance comparison of popular models across detailed metrics in **NR**, sub-category of PRMBENCH. The best performance for each metric is highlighted in **bold**, while the second-best performance is underlined. [†]: To reduce costs, we evaluated only a subset of 394 samples for o1-mini and Deepseek-R1.

Model Name	PRMScore	F1	Negative F1	Acc	Positive Acc	Negative Acc	First	similarity
<i>Open-source Process Level Reward Models</i>								
Skywork-PRM-1.5B	35.8	34.7	36.9	35.8	22.7	75.2	<u>71.4</u>	72.9
Skywork-PRM-7B	41.2	44.2	38.1	41.3	30.9	<u>72.6</u>	71.6	73.5
Llemma-PRM800k-7B	53.4	65.5	41.3	56.6	54.9	61.4	15.7	65.6
Llemma-MetaMath-7B	50.5	70.8	30.3	58.8	66.4	36.0	15.7	74.8
Llemma-oprm-7B	49.3	66.3	32.2	55.0	59.0	42.9	10.2	73.3
MATHMinos-Mistral-7B	54.0	70.8	37.2	60.2	66.6	43.2	40.6	79.1
MathShepherd-Mistral-7B	50.3	60.4	<u>40.2</u>	52.3	49.9	58.6	50.3	79.0
ReasonEval-7B	50.1	80.8	19.4	69.0	89.8	13.6	19.0	88.1
ReasonEval-34B	48.1	76.7	19.5	63.9	81.9	16.0	19.9	83.4
RLHFlow-PRM-Mistral-8B	47.3	79.8	14.9	67.3	88.7	10.5	18.9	93.7
RLHFlow-PRM-Deepseek-8B	48.9	82.3	15.4	70.7	93.6	9.8	15.6	94.6
Qwen2.5-Math-PRM-7B	55.1	84.5	25.7	<u>74.3</u>	96.1	16.3	22.7	91.6
Qwen2.5-Math-PRM-72B	58.8	<u>84.5</u>	33.1	74.8	<u>94.3</u>	22.8	33.5	89.4
Pure-PRM-7B	<u>55.2</u>	82.6	27.8	72.0	91.6	19.7	24.3	88.0
Avg.	49.9	70.3	29.4	60.9	70.5	35.6	30.7	81.9
<i>Open LLMs, Prompted as Critic Models</i>								
MetaMath-7B	46.9	74.9	19.0	61.7	75.2	18.7	13.6	77.8
MetaMath-13B	44.4	76.4	12.4	62.8	81.0	10.2	7.3	83.8
Qwen2.5-Math-72B	54.9	<u>81.2</u>	28.6	<u>70.2</u>	<u>89.5</u>	21.1	27.7	81.5
QwQ-Preview-32B	<u>55.6</u>	76.1	35.2	65.0	77.9	33.1	37.2	75.8
R1-Distill-Llama3.1-70B	48.1	84.4	11.8	73.5	98.1	6.6	8.5	95.6
R1-Distill-Qwen-7B	37.9	38.2	<u>37.7</u>	37.9	25.8	73.4	73.3	27.2
DeepSeek-R1 [†]	62.7	76.4	49.0	67.7	71.6	<u>57.3</u>	<u>61.7</u>	-
Avg.	50.1	72.5	27.7	62.7	74.2	31.5	32.8	73.6
<i>Proprietary LLMs, Prompted as Critic Models</i>								
GPT-4o	62.4	73.5	51.3	65.6	65.4	66.2	80.6	59.9
o1-mini [†]	<u>63.7</u>	80.4	<u>47.0</u>	71.4	80.6	<u>46.6</u>	47.9	-
Gemini-2.0-flash-exp	58.1	81.8	34.5	<u>71.5</u>	88.1	27.4	34.7	79.4
Gemini-2.0-thinking-exp-1219	63.8	<u>81.2</u>	46.4	72.2	<u>82.8</u>	44.0	<u>54.8</u>	74.6
Avg.	62.0	79.2	44.8	70.2	79.2	46.1	54.5	71.3

Table 17: A performance comparison of popular models across detailed metrics in **NCL**, sub-category of PRM-BENCH. The best performance for each metric is highlighted in **bold**, while the second-best performance is underlined. [†]: To reduce costs, we evaluated only a subset of 394 samples for o1-mini and Deepseek-R1.

Model Name	PRMScore	F1	Negative F1	Acc	Positive Acc	Negative Acc	First	similarity
<i>Open-source Process Level Reward Models</i>								
Skywork-PRM-1.5B	32.4	32.9	32.0	32.4	20.3	<u>85.4</u>	<u>76.4</u>	72.2
Skywork-PRM-7B	36.7	39.1	34.4	36.8	24.9	88.9	82.5	67.6
Llemma-PRM800k-7B	56.4	76.8	36.1	65.9	69.2	51.7	20.0	81.2
Llemma-MetaMath-7B	51.9	78.0	25.8	66.0	73.9	31.7	15.4	83.3
Llemma-oprm-7B	54.2	77.3	31.1	65.9	71.4	41.6	15.6	83.5
MATHMinos-Mistral-7B	57.0	77.3	36.7	66.6	71.1	48.6	39.7	82.0
MathShepherd-Mistral-7B	49.4	62.4	36.4	52.7	48.9	68.0	58.0	81.6
ReasonEval-7B	62.1	89.7	34.6	82.2	96.6	23.7	37.4	91.2
ReasonEval-34B	66.4	83.3	49.4	74.9	78.1	61.9	61.1	81.5
RLHFlow-PRM-Mistral-8B	56.6	85.8	27.4	76.2	89.6	22.5	24.5	90.5
RLHFlow-PRM-Deepseek-8B	55.7	87.8	23.5	79.0	94.6	16.2	21.4	93.5
Qwen2.5-Math-PRM-7B	<u>71.8</u>	90.8	<u>52.8</u>	<u>84.6</u>	<u>94.8</u>	43.3	56.2	83.7
Qwen2.5-Math-PRM-72B	73.7	<u>90.7</u>	56.8	84.7	93.2	50.6	66.9	81.4
Pure-PRM-7B	71.1	89.5	52.6	82.9	91.5	48.0	55.1	82.6
Avg.	56.8	75.8	37.8	67.9	72.7	48.7	45.0	82.6
<i>Open LLMs, Prompted as Critic Models</i>								
MetaMath-7B	47.3	89.8	4.7	81.6	97.3	2.7	1.3	97.2
MetaMath-13B	47.8	<u>90.0</u>	5.6	81.9	98.4	3.1	3.1	98.6
Qwen2.5-Math-72B	55.5	88.6	22.4	80.2	<u>99.1</u>	13.1	16.2	94.7
QwQ-Preview-32B	<u>67.4</u>	87.6	<u>47.2</u>	79.9	92.8	<u>38.1</u>	<u>42.5</u>	83.5
R1-Distill-Llama3.1-70B	61.4	89.8	33.0	<u>82.3</u>	97.1	22.1	26.8	90.9
R1-Distill-Qwen-7B	47.3	90.2	4.4	82.3	99.1	2.4	3.3	98.2
DeepSeek-R1 [†]	68.2	84.3	52.0	76.3	78.0	69.0	77.1	-
Avg.	56.4	88.6	24.2	80.6	94.5	21.5	24.3	93.9
<i>Proprietary LLMs, Prompted as Critic Models</i>								
GPT-4o	72.0	88.9	55.2	82.2	<u>88.8</u>	55.4	63.2	78.7
o1-mini [†]	74.5	88.9	60.0	<u>82.7</u>	85.6	69.8	<u>75.0</u>	-
Gemini-2.0-flash-exp	70.4	85.3	55.4	77.9	80.0	<u>69.4</u>	76.3	77.1
Gemini-2.0-thinking-exp-1219	<u>72.9</u>	89.4	<u>56.4</u>	83.0	89.8	55.5	64.8	79.5
Avg.	72.4	88.1	56.7	81.4	86.0	62.5	69.8	78.4

Table 18: A performance comparison of popular models across detailed metrics in **ES**, sub-category of PRMBENCH. The best performance for each metric is highlighted in **bold**, while the second-best performance is underlined. [†]: To reduce costs, we evaluated only a subset of 394 samples for o1-mini and Deepseek-R1.

Model Name	PRMScore	F1	Negative F1	Acc	Positive Acc	Negative Acc	First	similarity
<i>Open-source Process Level Reward Models</i>								
Skywork-PRM-1.5B	64.9	91.1	38.7	84.4	88.7	47.1	58.3	91.2
Skywork-PRM-7B	67.1	91.0	43.2	84.4	87.7	56.6	<u>69.6</u>	88.9
Llemma-PRM800k-7B	47.1	75.2	18.9	62.1	64.3	42.5	24.1	80.5
Llemma-MetaMath-7B	47.6	81.3	13.9	69.2	74.5	23.9	15.5	85.0
Llemma-oprm-7B	46.8	77.6	16.1	64.6	68.3	32.6	16.8	83.4
MATHMinos-Mistral-7B	52.1	80.7	23.5	69.2	72.6	42.4	39.4	85.5
MathShepherd-Mistral-7B	44.5	64.9	24.2	52.0	49.9	68.6	63.0	81.9
ReasonEval-7B	65.9	94.1	37.6	89.2	96.6	29.4	40.1	92.0
ReasonEval-34B	60.3	83.7	36.9	74.1	74.8	<u>68.2</u>	70.6	81.5
RLHFlow-PRM-Mistral-8B	55.1	90.0	20.2	82.2	90.0	20.1	22.9	92.3
RLHFlow-PRM-Deepseek-8B	55.0	92.4	17.7	86.1	<u>95.2</u>	13.4	16.6	95.1
Qwen2.5-Math-PRM-7B	67.3	<u>93.4</u>	41.2	<u>88.1</u>	94.4	37.6	43.8	87.9
Qwen2.5-Math-PRM-72B	71.1	92.9	49.2	87.5	91.6	54.5	61.6	83.1
Pure-PRM-7B	<u>68.8</u>	91.7	<u>45.8</u>	85.7	89.6	54.5	61.0	83.5
Avg.	58.1	85.7	30.5	77.1	81.3	42.2	43.1	86.6
<i>Open LLMs, Prompted as Critic Models</i>								
MetaMath-7B	48.9	92.6	5.1	86.3	94.7	4.0	4.1	96.3
MetaMath-13B	47.4	<u>93.6</u>	1.2	88.1	<u>97.8</u>	0.7	0.9	99.6
Qwen2.5-Math-72B	<u>71.6</u>	93.5	<u>49.6</u>	<u>88.5</u>	95.1	43.8	49.5	86.9
QwQ-Preview-32B	72.3	91.3	53.3	85.3	88.8	<u>62.4</u>	65.0	79.7
R1-Distill-Llama3.1-70B	65.5	93.4	37.6	88.1	95.3	31.6	37.2	89.6
R1-Distill-Qwen-7B	54.1	94.0	14.1	88.8	98.9	8.2	9.2	97.0
DeepSeek-R1 [†]	68.5	90.5	46.5	83.8	86.0	65.7	<u>63.2</u>	75.0
Avg.	61.2	92.7	29.6	87.0	93.8	30.9	32.7	89.2
<i>Proprietary LLMs, Prompted as Critic Models</i>								
GPT-4o	<u>69.7</u>	<u>89.9</u>	<u>49.6</u>	<u>83.1</u>	84.3	<u>74.3</u>	<u>76.9</u>	76.3
o1-mini [†]	67.7	89.7	45.7	82.7	<u>84.4</u>	68.6	70.8	74.2
Gemini-2.0-flash-exp	65.7	86.0	45.4	77.7	77.0	83.1	85.6	77.3
Gemini-2.0-thinking-exp-1219	71.3	91.4	51.2	85.4	87.6	68.5	72.8	81.0
Avg.	68.6	89.3	48.0	82.3	83.3	73.6	76.5	77.2

Table 19: A performance comparison of popular models across detailed metrics in SC, sub-category of PRMBENCH. The best performance for each metric is highlighted in **bold**, while the second-best performance is underlined. [†]: To reduce costs, we evaluated only a subset of 394 samples for o1-mini and Deepseek-R1.

Model Name	PRMScore	F1	Negative F1	Acc	Positive Acc	Negative Acc	First	similarity
<i>Open-source Process Level Reward Models</i>								
Skywork-PRM-1.5B	63.3	89.6	37.0	82.2	87.9	42.1	49.9	90.5
Skywork-PRM-7B	<u>67.7</u>	90.2	<u>45.2</u>	83.3	87.3	55.4	<u>65.4</u>	88.3
Llemma-PRM800k-7B	46.7	71.6	21.9	58.3	59.9	47.1	20.5	74.3
Llemma-MetaMath-7B	44.4	76.0	12.9	62.3	68.0	22.5	12.9	79.1
Llemma-oprm-7B	44.5	73.4	15.6	59.6	63.8	30.2	13.6	79.4
MATHMinos-Mistral-7B	50.7	76.1	25.3	63.8	66.6	45.8	39.0	80.6
MathShepherd-Mistral-7B	41.3	59.2	23.3	46.7	44.6	<u>60.5</u>	52.2	79.8
ReasonEval-7B	61.5	92.1	31.0	<u>85.8</u>	95.4	23.8	30.5	91.2
ReasonEval-34B	57.8	80.3	35.3	69.7	71.0	61.9	57.9	77.5
RLHFlow-PRM-Mistral-8B	54.4	87.8	21.0	78.9	87.8	20.9	18.8	90.4
RLHFlow-PRM-Deepseek-8B	53.2	90.4	15.9	82.8	<u>93.7</u>	12.2	11.6	94.2
Qwen2.5-Math-PRM-7B	66.3	91.5	41.2	85.1	92.2	39.1	42.0	85.0
Qwen2.5-Math-PRM-72B	72.2	<u>91.7</u>	52.7	85.9	90.1	58.7	65.5	80.5
Pure-PRM-7B	64.0	89.8	38.3	82.5	89.0	40.7	42.5	84.2
Avg.	56.3	82.8	29.8	73.4	78.4	40.1	37.3	83.9
<i>Open LLMs, Prompted as Critic Models</i>								
MetaMath-7B	48.4	90.9	5.8	83.5	96.5	3.6	3.8	96.7
MetaMath-13B	49.4	92.1	6.7	85.4	<u>98.1</u>	3.9	3.8	98.8
Qwen2.5-Math-72B	58.1	91.8	24.3	85.3	97.0	16.4	15.0	92.6
QwQ-Preview-32B	<u>66.2</u>	87.0	<u>45.5</u>	78.9	82.9	<u>57.0</u>	<u>51.6</u>	74.4
R1-Distill-Llama3.1-70B	65.8	91.9	39.7	85.7	93.4	35.4	35.5	86.1
R1-Distill-Qwen-7B	48.4	<u>92.1</u>	4.7	85.4	98.4	2.7	2.8	98.4
DeepSeek-R1 [†]	73.5	91.3	55.6	<u>85.5</u>	85.7	84.1	77.1	75.4
Avg.	58.5	91.0	26.0	84.2	93.1	29.0	27.1	88.9
<i>Proprietary LLMs, Prompted as Critic Models</i>								
GPT-4o	70.7	88.2	<u>53.3</u>	81.2	81.3	<u>80.4</u>	<u>81.4</u>	71.4
o1-mini [†]	73.8	92.2	55.5	86.7	87.7	77.8	74.5	80.0
Gemini-2.0-flash-exp	66.0	83.6	48.4	75.1	73.2	87.5	87.3	70.8
Gemini-2.0-thinking-exp-1219	<u>71.0</u>	<u>88.8</u>	53.2	<u>81.9</u>	<u>82.6</u>	77.2	79.2	73.7
Avg.	70.4	88.2	52.6	81.2	81.2	80.7	80.6	74.0

Table 20: A performance comparison of popular models across detailed metrics in **DC**, sub-category of PRMBENCH. The best performance for each metric is highlighted in **bold**, while the second-best performance is underlined. [†]: To reduce costs, we evaluated only a subset of 394 samples for the o1-mini and Deepseek-R1.

Model Name	PRMScore	F1	Negative F1	Acc	Positive Acc	Negative Acc	First	similarity
<i>Open-source Process Level Reward Models</i>								
Skywork-PRM-1.5B	66.5	91.3	41.7	84.8	89.7	47.3	62.2	92.5
Skywork-PRM-7B	69.9	91.5	48.2	85.4	88.8	59.3	71.5	90.8
Llemma-PRM800k-7B	53.3	79.9	26.6	68.5	70.9	50.0	29.4	86.2
Llemma-MetaMath-7B	52.1	83.7	20.4	73.0	78.5	30.3	15.6	88.9
Llemma-oprm-7B	53.5	81.7	25.4	70.6	74.0	43.8	27.2	89.0
MATHMinos-Mistral-7B	57.8	82.5	33.1	72.3	74.5	56.3	53.8	87.3
MathShepherd-Mistral-7B	47.7	66.6	28.7	54.5	51.7	75.1	72.2	87.3
ReasonEval-7B	66.0	93.7	38.3	88.6	96.9	29.0	40.6	92.8
ReasonEval-34B	67.5	87.7	47.2	80.1	81.0	<u>73.3</u>	<u>73.3</u>	86.9
RLHFlow-PRM-Mistral-8B	63.8	91.0	36.7	84.2	90.7	37.5	45.0	90.9
RLHFlow-PRM-Deepseek-8B	66.2	93.0	39.5	87.5	94.9	33.5	43.3	92.1
Qwen2.5-Math-PRM-7B	<u>78.5</u>	94.9	<u>62.0</u>	91.0	<u>95.2</u>	60.5	70.4	85.4
Qwen2.5-Math-PRM-72B	78.6	<u>94.7</u>	62.5	<u>90.7</u>	94.3	64.1	73.7	84.8
Pure-PRM-7B	76.9	93.9	59.9	89.5	92.9	64.7	70.8	86.3
Avg.	64.2	87.6	40.7	80.0	83.9	51.8	53.5	88.7
<i>Open LLMs, Prompted as Critic Models</i>								
MetaMath-7B	48.8	93.9	3.7	88.5	97.1	2.4	2.1	97.7
MetaMath-13B	48.1	94.7	1.5	89.9	99.1	0.8	0.6	99.0
Qwen2.5-Math-72B	59.1	93.1	25.1	87.4	<u>99.6</u>	14.7	19.8	96.2
QwQ-Preview-32B	<u>66.9</u>	91.7	<u>42.2</u>	85.4	94.8	<u>34.2</u>	<u>43.1</u>	89.1
R1-Distill-Llama3.1-70B	61.1	93.9	28.2	88.8	98.3	18.5	22.0	94.8
R1-Distill-Qwen-7B	48.0	<u>94.6</u>	1.3	<u>89.8</u>	99.8	0.7	0.7	99.5
DeepSeek-R1 [†]	75.4	92.7	58.1	87.6	88.8	78.1	82.4	82.7
Avg.	58.2	93.5	22.9	88.2	96.8	21.4	24.4	94.2
<i>Proprietary LLMs, Prompted as Critic Models</i>								
GPT-4o	71.1	<u>92.1</u>	50.2	<u>86.3</u>	<u>90.4</u>	56.5	60.4	85.2
o1-mini [†]	72.3	91.6	53.1	85.7	87.3	72.7	<u>75.0</u>	80.1
Gemini-2.0-flash-exp	67.3	88.3	46.4	80.8	82.4	<u>68.8</u>	76.0	84.9
Gemini-2.0-thinking-exp-1219	<u>71.8</u>	92.6	<u>51.1</u>	87.1	91.4	55.8	60.5	86.1
Avg.	70.7	91.1	50.2	85.0	87.9	63.5	68.0	84.1

Table 21: A performance comparison of popular models across detailed metrics in **CI** sub-category of PRMBENCH. The best performance for each metric is highlighted in **bold**, while the second-best performance is underlined. [†]: To reduce costs, we evaluated only a subset of 394 samples for o1-mini and Deepseek-R1.

Model Name	PRMScore	F1	Negative F1	Acc	Positive Acc	Negative Acc	First	similarity
<i>Open-source Process Level Reward Models</i>								
Skywork-PRM-1.5B	33.1	36.2	30.1	33.3	23.2	78.4	73.9	70.2
Skywork-PRM-7B	36.8	42.9	30.8	37.4	28.8	<u>76.0</u>	<u>72.8</u>	70.2
Llemma-PRM800k-7B	51.0	71.7	30.3	59.7	62.3	47.9	24.4	75.0
Llemma-MetaMath-7B	50.5	78.7	22.4	66.5	75.5	26.5	15.6	82.8
Llemma-oprm-7B	49.2	74.4	24.1	61.7	68.0	33.4	16.1	81.2
MATHMinos-Mistral-7B	52.8	77.5	28.0	65.7	73.3	34.2	29.5	83.1
MathShepherd-Mistral-7B	47.2	63.6	30.9	52.3	51.8	54.5	45.0	81.4
ReasonEval-7B	55.6	88.4	22.9	79.8	<u>95.4</u>	15.4	17.3	92.8
ReasonEval-34B	57.7	80.5	35.0	70.0	76.8	41.6	35.2	82.6
RLHFlow-PRM-Mistral-8B	51.5	85.3	17.6	75.0	89.9	13.7	16.2	93.1
RLHFlow-PRM-Deepseek-8B	49.0	87.0	10.9	77.4	94.4	7.1	6.7	95.9
Qwen2.5-Math-PRM-7B	57.6	88.8	26.5	80.5	95.6	18.1	17.7	91.7
Qwen2.5-Math-PRM-72B	60.3	<u>88.5</u>	32.2	<u>80.3</u>	93.8	24.1	24.9	89.4
Pure-PRM-7B	<u>60.3</u>	87.4	<u>33.1</u>	78.8	91.4	27.0	24.0	87.3
Avg.	50.9	75.0	26.8	65.6	72.9	35.6	29.9	84.0
<i>Open LLMs, Prompted as Critic Models</i>								
MetaMath-7B	46.5	88.0	5.0	78.7	96.3	3.0	2.3	96.2
MetaMath-13B	49.0	88.6	9.4	79.7	97.6	5.4	3.1	96.7
Qwen2.5-Math-72B	47.4	88.5	6.3	79.5	99.9	3.3	3.1	98.8
QwQ-Preview-32B	<u>57.8</u>	86.4	<u>29.1</u>	77.2	92.8	<u>21.5</u>	<u>18.3</u>	88.3
R1-Distill-Llama3.1-70B	48.8	<u>88.9</u>	8.7	<u>80.3</u>	98.7	4.8	4.8	97.2
R1-Distill-Qwen-7B	45.6	89.5	1.6	81.0	<u>99.3</u>	0.8	1.1	99.1
DeepSeek-R1 [†]	63.1	82.7	43.5	73.5	82.6	43.7	34.8	-
Avg.	51.2	87.5	14.8	78.6	95.3	11.8	9.6	96.1
<i>Proprietary LLMs, Prompted as Critic Models</i>								
GPT-4o	62.5	<u>86.8</u>	38.3	<u>78.2</u>	<u>88.8</u>	<u>34.7</u>	<u>33.3</u>	83.8
o1-mini [†]	<u>61.8</u>	84.9	<u>38.7</u>	75.7	88.6	33.1	18.8	-
Gemini-2.0-flash-exp	61.8	83.2	40.4	73.8	80.5	45.7	44.2	81.4
Gemini-2.0-thinking-exp-1219	60.3	87.1	33.5	78.4	90.5	28.0	27.0	87.0
Avg.	61.6	85.5	37.7	76.5	87.1	35.4	30.8	84.0

Table 22: A performance comparison of popular models across detailed metrics in **PS**, sub-category of PRMBENCH. The best performance for each metric is highlighted in **bold**, while the second-best performance is underlined. [†]: To reduce costs, we evaluated only a subset of 394 samples for o1-mini and Deepseek-R1.

Model Name	PRMScore	F1	Negative F1	Acc	Positive Acc	Negative Acc	First	similarity
<i>Open-source Process Level Reward Models</i>								
Skywork-PRM-1.5B	32.3	36.3	28.2	32.5	22.8	84.0	79.1	74.5
Skywork-PRM-7B	37.4	44.9	29.9	38.3	29.8	<u>83.3</u>	<u>77.8</u>	73.3
Llemma-PRM800k-7B	53.5	77.8	29.1	66.2	70.4	43.8	23.1	85.1
Llemma-MetaMath-7B	51.3	80.6	22.1	68.9	76.6	27.9	15.1	85.6
Llemma-oprm-7B	51.3	78.4	24.1	66.4	72.5	33.8	17.2	86.6
MATHMinos-Mistral-7B	55.8	79.1	32.4	68.1	72.8	45.2	41.4	83.7
MathShepherd-Mistral-7B	48.6	65.6	31.7	54.2	52.5	62.8	57.1	84.0
ReasonEval-7B	58.0	90.6	25.4	83.2	96.7	16.9	24.6	93.8
ReasonEval-34B	64.3	84.4	44.3	75.6	79.4	57.3	56.8	83.6
RLHFlow-PRM-Mistral-8B	56.2	87.5	24.9	78.6	90.3	21.0	27.1	92.2
RLHFlow-PRM-Deepseek-8B	55.4	89.5	21.4	81.5	95.1	14.9	19.4	95.0
Qwen2.5-Math-PRM-7B	69.1	91.7	46.6	85.6	<u>95.4</u>	37.3	46.9	87.0
Qwen2.5-Math-PRM-72B	71.2	<u>91.5</u>	50.9	<u>85.5</u>	93.9	44.3	55.9	85.2
Pure-PRM-7B	<u>69.2</u>	90.2	<u>48.3</u>	83.6	91.3	45.3	51.7	84.5
Avg.	55.3	77.7	32.8	69.2	74.3	44.1	42.4	85.3
<i>Open LLMs, Prompted as Critic Models</i>								
MetaMath-7B	48.3	90.4	6.2	82.5	95.3	4.1	3.5	96.6
MetaMath-13B	48.1	<u>91.4</u>	4.8	<u>84.3</u>	98.7	2.6	4.0	99.6
Qwen2.5-Math-72B	53.8	90.5	17.1	82.9	<u>99.3</u>	9.6	11.5	96.3
QwQ-Preview-32B	<u>62.7</u>	89.2	<u>36.1</u>	81.6	94.9	<u>26.6</u>	<u>29.8</u>	88.8
R1-Distill-Llama3.1-70B	54.1	91.0	17.2	83.7	99.0	9.9	12.0	96.1
R1-Distill-Qwen-7B	46.8	91.8	1.9	84.9	99.5	1.0	1.2	99.3
DeepSeek-R1 [†]	69.2	89.0	49.5	81.9	87.3	54.1	59.5	-
Avg.	54.7	90.5	19.0	83.1	96.3	15.4	17.4	96.1
<i>Proprietary LLMs, Prompted as Critic Models</i>								
GPT-4o	65.7	<u>89.2</u>	42.2	<u>81.8</u>	<u>90.5</u>	39.3	41.3	84.8
o1-mini [†]	64.8	86.7	<u>42.9</u>	78.4	84.5	<u>48.2</u>	<u>43.8</u>	-
Gemini-2.0-flash-exp	66.2	86.3	46.1	78.1	82.7	55.5	60.2	83.4
Gemini-2.0-thinking-exp-1219	<u>65.7</u>	89.7	41.8	82.5	91.8	37.0	40.2	86.4
Avg.	65.6	88.0	43.2	80.2	87.4	45.0	46.4	84.9

Table 23: A performance comparison of popular models across detailed metrics in **DR**. sub-category of PRMBENCH. The best performance for each metric is highlighted in **bold**, while the second-best performance is underlined. [†]: To reduce costs, we evaluated only a subset of 394 samples for o1-mini and Deepseek-R1.

1. Original Process

Question: In 1992, a scoop of gelato could be purchased in Italy for 1200 lire. The same gelato would have cost \$1.50 in the U.S. At the equivalent exchange rate between the lire and the dollar, how many dollars would be equivalent to 1,000,000 lire?

Step 1: Let's call the number of dollars x . So, the problem is asking us to solve the equation $x = 1000000 * \frac{1.50}{1200}$.

Step 2: There are two ways to solve this equation. We could either divide 1000000 by 1200 and then multiply the result by 1.50 or we could divide 1.50 by 1200 and then multiply the result by 1000000.

Step 3: I think the second way will be easier. $\frac{1.50}{1200} = 0.00125$.

Step 4: That's true. Now we just have to solve the equation $x = 1000000 * 0.00125$.

Step 5: So $x = 1250$.

Step 6: That's the final answer.

Step 7: Right. So 1,000,000 lire is equivalent to \$1250.

2. Modified Process

Question: In 1992, a scoop of gelato could be purchased in Italy for 1200 lire. The same gelato would have cost \$1.50 in the U.S. At the equivalent exchange rate between the lire and the dollar, how many dollars would be equivalent to 1,000,000 lire?

Step 1: Let's call the number of dollars x . So, the problem is asking us to solve the equation $x = 1000000 * \frac{1.50}{1200}$.

Step 2: There are two ways to solve this equation. We could either divide 1000000 by 1200 and then multiply the result by 1.50 or we could divide 1.50 by 1200 and then multiply the result by 1000000.

Step 3: I think the second way will be easier. $\frac{1.50}{1200} = 0.00125$.

Step 4: That's true. But let's confirm the computation again: $\frac{1.50}{1200} = 0.00125$.

Step 5: Now we just have to solve the equation $x = 1000000 * 0.00125$.

Step 6: So $x = 1250$.

Step 7: Let's double-check by performing the multiplication again: $1000000 * 0.00125 = 1250$.

Step 8: That's the final answer.

Step 9: Right. So 1,000,000 lire is equivalent to \$1250.

3. Reason

Step 4 unnecessarily repeats the computation of $\frac{1.50}{1200}$, which was already completed in Step 3. Step 7 redundantly re-checks the multiplication, adding no new information. These redundant steps make the reasoning process less concise without affecting correctness.

Figure 7: An example data instance for **NR** category.

1. Original Process

Question: If $x^2 - x - 1 = 0$, what is the value of $x^3 - 2x + 1$?

Step 1: Let's start with the first equation and see if we can solve for x.

Step 2: Ok. We can use the quadratic formula to find x.

Step 3: The quadratic formula is $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$.

Step 4: Yes, that's correct. In this case a , b , and c are 1, -1, and -1, respectively.

Step 5: So if we plug those values in, we get $x = \frac{1 \pm \sqrt{1^2 - 4 \cdot 1 \cdot -1}}{2 \cdot 1}$.

Step 6: That simplifies to $x = \frac{1 \pm \sqrt{5}}{2}$. So x can be either $\frac{1 + \sqrt{5}}{2}$ or $\frac{1 - \sqrt{5}}{2}$.

Step 7: Now let's move to the second equation.

Step 8: Yes, we need to find $x^3 - 2x + 1$.

Step 9: Right. Well since we now know that $x^2 - x - 1 = 0$, we can replace x^2 in the second equation with $x + 1$.

Step 10: That means that $x^3 - 2x + 1 = x(x + 1) - 2x + 1$.

Step 11: And if we substitute in the possible values of x that we found earlier, we get either $\frac{1 + \sqrt{5}}{2}(\frac{1 + \sqrt{5}}{2} + 1) - 2\frac{1 + \sqrt{5}}{2} + 1$ or $\frac{1 - \sqrt{5}}{2}(\frac{1 - \sqrt{5}}{2} + 1) - 2\frac{1 - \sqrt{5}}{2} + 1$.

Step 12: Let's start with the first case: $\frac{1 + \sqrt{5}}{2}(\frac{1 + \sqrt{5}}{2} + 1) - 2\frac{1 + \sqrt{5}}{2} + 1$.

Step 13: There's a shorter way. We have $x^3 - 2x + 1 = x(x + 1) - 2x + 1$.

Step 14: $x^3 - 2x + 1 = (x^2 + x) - 2x + 1$.

Step 15: And we know that $x^2 - x - 1 = 0$, so $x^2 = x + 1$.

Step 16: So we can just substitute $x + 1$ for x^2 .

Step 17: So $x^3 - 2x + 1 = (x^2 + x) - 2x + 1$ becomes $x^3 - 2x + 1 = (x + 1 + x) - 2x + 1$.

Step 18: And that simplifies to $x^3 - 2x + 1 = (2x + 1) - 2x + 1$.

Step 19: So $x^3 - 2x + 1 = 2$.

Step 20: So, we don't need to do the second case because we already have the answer.

Step 21: Right. So if $x^2 - x - 1 = 0$ then $x^3 - 2x + 1 = 2$.

2. Modified Process

Question: If $x^2 - x - 1 = 0$, what is the value of $x^3 - 2x + 1$?

Step 1: Let's start with the first equation and see if we can solve for x.

Step 2: Ok. We can use the quadratic formula to find x.

Step 3: The quadratic formula is $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$.

Step 4: Yes, that's correct. In this case a , b , and c are 1, -1, and -1, respectively.

Step 5: So if we plug those values in, we get $x = \frac{1 \pm \sqrt{1^2 - 4 \cdot 1 \cdot -1}}{2 \cdot 1}$.

Step 6: That simplifies to $x = \frac{1 \pm \sqrt{5}}{2}$. So x can be either $\frac{1 + \sqrt{5}}{2}$ or $\frac{1 - \sqrt{5}}{2}$.

Step 7: Now let's move to the second equation.

Step 8: Assume $x^3 - 2x + 1 = 2$ for verification.

Step 9: Substituting $x^3 - 2x + 1 = 2$ into the equation, it implies $x(x^2) - 2x + 1 = 2$.

Step 10: Since $x^2 = x + 1$, $x(x + 1) - 2x + 1 = 2$.

Step 11: Simplifying confirms $x^3 - 2x + 1 = 2$, verifying the assumption.

Step 12: Therefore, $x^3 - 2x + 1 = 2$ is correct.

3. Reason

Steps 8 and 9 assume the result ($x^3 - 2x + 1 = 2$) before deriving it, and step 11 uses circular logic to verify the assumption by substituting it back into the equation. This process does not independently solve the problem and creates a flawed reasoning chain.

Figure 8: An example data instance for NCL category.

1. Original Process

Question: In 1992, a scoop of gelato could be purchased in Italy for 1200 lire. The same gelato would have cost \$1.50 in the U.S. At the equivalent exchange rate between the lire and the dollar, how many dollars would be equivalent to 1,000,000 lire?

Step 1: Let's call the number of dollars x . So, the problem is asking us to solve the equation $x = 1000000 * \frac{1.50}{1200}$.

Step 2: There are two ways to solve this equation. We could either divide 1000000 by 1200 and then multiply the result by 1.50 or we could divide 1.50 by 1200 and then multiply the result by 1000000.

Step 3: I think the second way will be easier. $\frac{1.50}{1200} = 0.00125$.

Step 4: That's true. Now we just have to solve the equation $x = 1000000 * 0.00125$.

Step 5: So $x=1250$.

Step 6: That's the final answer.

Step 7: Right. So 1,000,000 lire is equivalent to \$1250.

2. Modified Process

Question: In 1992, a scoop of gelato could be purchased in Italy for 1200 lire. The same gelato would have cost \$1.50 in the U.S. At the equivalent exchange rate between the lire and the dollar, how many dollars would be equivalent to 1,000,000 lire?

Step 1: Let's call the number of dollars x . So, the problem is asking us to solve the equation $x = 1000000 * \frac{1.50}{1200}$.

Step 2: There are two ways to solve this equation. We could either divide 1000000 by 1200 and then multiply the result by 1.50 or we could divide 1.50 by 1200 and then multiply the result by 1000000.

Step 3: I think the second way will be easier. $\frac{1.50}{1200} = 0.00125$.

Step 4: That's true. Now we just have to solve the equation $x = 1000000 * 0.00125$.

Step 5: So $x=1250$.

Step 6: But wait, we need to account for transaction fees that add an additional 10%.

Step 7: So $x = 1250 + 0.10 * 1250$.

Step 8: That gives $x=1375$.

Step 9: Right. So 1,000,000 lire is equivalent to \$1375.

3. Reason

A counterfactual step was introduced in step 6 by claiming the need to account for an additional 10% transaction fee, which is not mentioned in the problem. This erroneous assumption leads to the incorrect result of \$1375 instead of \$1250.

Figure 9: An example data instance for ES. category

1. Original Process

Question: If $x^2 - x - 1 = 0$, what is the value of $x^3 - 2x + 1$?

Step 1: Let's start with the first equation and see if we can solve for x .

Step 2: We can use the quadratic formula to find x .

Step 3: The quadratic formula is $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$.

Step 4: Substituting $a = 1$, $b = -1$, and $c = -1$, we get $x = \frac{1 \pm \sqrt{5}}{2}$.

Step 5: Since $x^2 - x - 1 = 0$, we replace x^2 with $x + 1$ in $x^3 - 2x + 1$.

Step 6: Substituting, $x^3 - 2x + 1 = x(x + 1) - 2x + 1$.

Step 7: Simplify: $x^3 - 2x + 1 = (x^2 + x) - 2x + 1$.

Step 8: Substitute $x^2 = x + 1$: $x^3 - 2x + 1 = (x + 1 + x) - 2x + 1$.

Step 9: Simplify: $x^3 - 2x + 1 = 2$.

2. Modified Process

Question: If $x^2 - x - 1 = 0$, what is the value of $x^3 - 2x + 1$?

Step 1: Let's start with the first equation and see if we can solve for x .

Step 2: We can use the quadratic formula to find x .

Step 3: The quadratic formula is $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$.

Step 4: Substituting $a = 1$, $b = -1$, and $c = -1$, we get $x = \frac{1 \pm \sqrt{5}}{2}$.

Step 5: Since $x^2 - x - 1 = 0$, we replace x^2 with $x + 1$ in $x^3 - 2x + 1$.

Step 6: Substituting, $x^3 - 2x + 1 = x(x + 1) - 2x + 1$.

Step 7: Simplify: $x^3 - 2x + 1 = (x^2 + x) - 2x + 1$.

Step 8: Incorrectly substitute $x^2 = x - 1$ instead of $x + 1$: $x^3 - 2x + 1 = (x - 1 + x) - 2x + 1$.

Step 9: Simplify: $x^3 - 2x + 1 = 0$.

3. Reason

Step 8 introduces a step contradiction by incorrectly substituting $x^2 = x - 1$ instead of the correct $x^2 = x + 1$. This conflicts with the earlier derivation from $x^2 - x - 1 = 0$, where $x^2 = x + 1$. As a result, the final answer becomes 0 instead of the correct value, 2.

Figure 10: An example data instance for SC. category.

1. Original Process

Question: If $x^2 - x - 1 = 0$, what is the value of $x^3 - 2x + 1$?

Step 1: Let's start with the first equation and see if we can solve for x.

Step 2: We can use the quadratic formula to find x.

Step 3: The quadratic formula is $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$.

Step 4: Substituting $a = 1$, $b = -1$, and $c = -1$, we get $x = \frac{1 \pm \sqrt{5}}{2}$.

Step 5: Since $x^2 - x - 1 = 0$, we replace x^2 with $x + 1$ in $x^3 - 2x + 1$.

Step 6: Substituting, $x^3 - 2x + 1 = x(x + 1) - 2x + 1$.

Step 7: Simplify: $x^3 - 2x + 1 = (x^2 + x) - 2x + 1$.

Step 8: Substitute $x^2 = x + 1$: $x^3 - 2x + 1 = (x + 1 + x) - 2x + 1$.

Step 9: Simplify: $x^3 - 2x + 1 = 2$.

2. Modified Process

Question: If $x^2 - x - 1 = 0$, what is the value of $x^3 - 2x + 1$?

Step 1: Let's start with the first equation and see if we can solve for x.

Step 2: We can use the quadratic formula to find x.

Step 3: The quadratic formula is $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$.

Step 4: Substituting $a = 1$, $b = -1$, and $c = -1$, we get $x = \frac{1 \pm \sqrt{5}}{2}$.

Step 5: Since $x^2 - x - 1 = 0$, we replace x^2 with $x + 1$ in $x^3 - 2x + 1$.

Step 6: Substituting, $x^3 - 2x + 1 = x(x + 1) - 2x + 1$.

Step 7: To simplify x^3 , use the property of cubic roots: $x^3 = \sqrt[3]{x}$.

Step 8: Substitute $x^3 = \sqrt[3]{x}$ and $x^2 = x + 1$ into $x^3 - 2x + 1$: $x^3 - 2x + 1 = \sqrt[3]{x} - 2x + 1$.

Step 9: Simplify further: $x^3 - 2x + 1 = \sqrt[3]{x} - x + 1$.

3. Reason

A domain inconsistency was introduced in Step 7 by incorrectly assuming that $x^3 = \sqrt[3]{x}$. This assumption is valid for root-based equations but not in the context of solving a polynomial equation. This error propagates through the rest of the reasoning, leading to a flawed solution.

Figure 11: An example data instance for **DC** category.

1. Original Process

Question: In 1992, a scoop of gelato could be purchased in Italy for 1200 lire. The same gelato would have cost \$1.50 in the U.S. At the equivalent exchange rate between the lire and the dollar, how many dollars would be equivalent to 1,000,000 lire?

Step 1: Let's call the number of dollars x . So, the problem is asking us to solve the equation $x = 1000000 * \frac{1.50}{1200}$.

Step 2: There are two ways to solve this equation. We could either divide 1000000 by 1200 and then multiply the result by 1.50 or we could divide 1.50 by 1200 and then multiply the result by 1000000.

Step 3: I think the second way will be easier. $\frac{1.50}{1200} = 0.00125$.

Step 4: That's true. Now we just have to solve the equation $x = 1000000 * 0.00125$.

Step 5: So $x = 1250$.

Step 6: That's the final answer.

Step 7: Right. So 1,000,000 lire is equivalent to \$1250.

2. Modified Process

Question: In 1992, a scoop of gelato could be purchased in Italy for 1200 lire. The same gelato would have cost \$1.50 in the U.S. At the equivalent exchange rate between the lire and the dollar, how many dollars would be equivalent to 1,000,000 lire?

Step 1: Let's call the number of dollars x . So, the problem is asking us to solve the equation $x = 1000000 * \frac{1.50}{1200}$.

Step 2: There are two ways to solve this equation. We could either divide 1000000 by 1200 and then multiply the result by 1.50 or we could divide 1.50 by 1200 and then multiply the result by 1000000.

Step 3: I think the second way will be easier. $\frac{1.50}{1200} = 0.00125$.

Step 4: That's true. Now we just have to solve the equation $x = 1000000 * 0.00125$.

Step 5: So $x = 1250$.

Step 6: That's the final answer.

Step 7: Actually, after a closer look, $x = 0.00125 \times 1000000$ should be 12500. I must have miscalculated earlier.

Step 8: So 1,000,000 lire is equivalent to \$12500.

3. Reason

Step 7 introduces a confident hallucination. It contradicts the correct value of 1250 by recalculating and asserting that the result is 12500, which is false. This error is presented with unwarranted certainty to simulate confident misjudgment.

Figure 12: An example data instance for **CI** category.

1. Original Process

Question: Find all values of x such that $|x - 3| = 5$.

Step 1: Start with the equation $|x - 3| = 5$.

Step 2: Recognize that the absolute value equation $|a| = b$ implies two cases: $a = b$ or $a = -b$.

Step 3: Apply this rule to $|x - 3| = 5$, resulting in two equations: $x - 3 = 5$ or $x - 3 = -5$.

Step 4: Solve each equation: $x = 8$ from $x - 3 = 5$ and $x = -2$ from $x - 3 = -5$.

Step 5: The solutions are $x = 8$ and $x = -2$.

2. Modified Process

Question: Find all values of x such that $(\sqrt{x - 3})^2 = 5$.

Step 1: $(\sqrt{x - 3})^2 = 5$ can be represented as a special form of $|x - 3| = 5$.

Step 2: Start with the equation $|x - 3| = 5$.

Step 3: Recognize that the absolute value equation $|a| = b$ implies two cases: $a = b$ or $a = -b$.

Step 4: Apply this rule to $|x - 3| = 5$, resulting in two equations: $x - 3 = 5$ or $x - 3 = -5$.

Step 5: Solve each equation: $x = 8$ from $x - 3 = 5$ and $x = -2$ from $x - 3 = -5$.

Step 6: The solutions are $x = 8$ and $x = -2$.

3. Reason

The first step introduces a missing condition by changing the equation to $|x - 3| = 5$, which ignores the implicit requirement that $x - 3$ must be non-negative. The final step incorrectly asserts that the solutions are $x = 8$ and $x = -2$ without considering the domain restriction.

Figure 13: An example data instance for **PS.** category.

1. Original Process

Question: Given that $\mathbf{F}(r) = (f(r), g(r), h(r))$ is a three-dimensional vector, prove that: $\forall r_1, r_2, \exists r^* \in (r_1, r_2)$ such that $\mathbf{F}(r_1) - \mathbf{F}(r_2) = (r_1 - r_2)\nabla\mathbf{F}(r^*)$.

Step 1: Understanding the Given Information and Problem We have a vector function $\mathbf{F}(r) = (f(r), g(r), h(r))$, where f, g , and h are real-valued functions of the real variable r . We need to prove that for any r_1, r_2 , there exists an r^* such that $r_1 < r^* < r_2$ and

$$\mathbf{F}(r_1) - \mathbf{F}(r_2) = (r_1 - r_2)\nabla\mathbf{F}(r^*),$$

where $\nabla\mathbf{F}(r^*)$ denotes the gradient of \mathbf{F} at r^* (i.e., a vector composed of the derivatives of each component with respect to r).

Step 2: Breaking Down the Problem into Components Express the vector equation in terms of its components:

$$(f(r_1) - f(r_2), g(r_1) - g(r_2), h(r_1) - h(r_2)) = (r_1 - r_2)(f'(r^*), g'(r^*), h'(r^*)).$$

Thus, we need to find an $r^* \in (r_1, r_2)$ such that for each component:

$$\frac{f(r_1) - f(r_2)}{r_1 - r_2} = f'(r^*), \quad \frac{g(r_1) - g(r_2)}{r_1 - r_2} = g'(r^*), \quad \frac{h(r_1) - h(r_2)}{r_1 - r_2} = h'(r^*).$$

Step 3: Considering the Application of the One-Dimensional Mean Value Theorem In one dimension, the Mean Value Theorem (MVT) states that for a function f continuous on $[a, b]$ and differentiable on (a, b) , there exists $c \in (a, b)$ such that:

$$f(b) - f(a) = f'(c)(b - a).$$

Applying the MVT to each component function f, g , and h , we obtain corresponding points c_f, c_g , and c_h . However, these points generally differ ($c_f \neq c_g \neq c_h$), so we cannot guarantee the existence of a common r^* where all component equations hold simultaneously.

Step 4: Considering a Counterexample Since we cannot guarantee a common r^* satisfying all components, let's attempt to find a specific counterexample to demonstrate that the original proposition may not hold.

Step 5: Trying the Function $\mathbf{F}(r) = (\sin r, \cos r, r)$ Let $\mathbf{F}(r) = (\sin r, \cos r, r)$. Choose two distinct real numbers, for example, $r_1 = 0$ and $r_2 = \pi$.

Step 6: Compute $\mathbf{F}(r_1) - \mathbf{F}(r_2)$:

$$\mathbf{F}(0) - \mathbf{F}(\pi) = (\sin 0 - \sin \pi, \cos 0 - \cos \pi, 0 - \pi) = (0 - 0, 1 - (-1), -\pi) = (0, 2, -\pi).$$

Step 7: Compute $(r_1 - r_2)\nabla\mathbf{F}(r^*)$:

$$(r_1 - r_2)(\cos r^*, -\sin r^*, 1) = (0 - \pi)(\cos r^*, -\sin r^*, 1) = (-\pi)(\cos r^*, -\sin r^*, 1).$$

Step 8: We need to find $r^* \in (0, \pi)$ such that:

$$(0, 2, -\pi) = (-\pi \cos r^*, \pi \sin r^*, -\pi).$$

Step 9: This leads to the equations:

$$-\pi \cos r^* = 0, \quad \pi \sin r^* = 2, \quad -\pi = -\pi.$$

Step 10: The third equation is always true. Simplifying the first two equations: From $-\pi \cos r^* = 0$: $\cos r^* = 0$.

Figure 14: An example data instance for **DR**. category (1).

Step 11: From $\pi \sin r^* = 2$:

$$\sin r^* = \frac{2}{\pi}.$$

Step 12: However, $\cos r^* = 0$ implies $r^* = \frac{\pi}{2}$. Substituting $r^* = \frac{\pi}{2}$ into the second equation:

$$\sin\left(\frac{\pi}{2}\right) = 1 \neq \frac{2}{\pi}.$$

Step 13: Since $\frac{2}{\pi} \approx 0.6366$, the equality does not hold. Therefore, there is no r^* in $(0, \pi)$ that satisfies both equations simultaneously.

Step 14: The above calculations show that for the function $\mathbf{F}(r) = (\sin r, \cos r, r)$, there does not exist an $r^* \in (r_1, r_2)$ such that:

$$\mathbf{F}(r_1) - \mathbf{F}(r_2) = (r_1 - r_2)\nabla\mathbf{F}(r^*).$$

2. Modified Process

Question: Given that $\mathbf{F}(r) = (f(r), g(r), h(r))$ is a three-dimensional vector, prove that: $\forall r_1, r_2, \exists r^* \in (r_1, r_2)$ such that $\mathbf{F}(r_1) - \mathbf{F}(r_2) = (r_1 - r_2)\nabla\mathbf{F}(r^*)$.

Step 1: Understanding the Given Information and Problem We have a vector function $\mathbf{F}(r) = (f(r), g(r), h(r))$, where f, g , and h are real-valued functions of the real variable r . We need to prove that for any r_1, r_2 , there exists an r^* such that $r_1 < r^* < r_2$ and

$$\mathbf{F}(r_1) - \mathbf{F}(r_2) = (r_1 - r_2)\nabla\mathbf{F}(r^*),$$

where $\nabla\mathbf{F}(r^*)$ denotes the gradient of \mathbf{F} at r^* (i.e., a vector composed of the derivatives of each component with respect to r).

Step 2: Breaking Down the Problem into Components Express the vector equation in terms of its components:

$$(f(r_1) - f(r_2), g(r_1) - g(r_2), h(r_1) - h(r_2)) = (r_1 - r_2)(f'(r^*), g'(r^*), h'(r^*)).$$

Thus, we need to find an $r^* \in (r_1, r_2)$ such that for each component:

$$\frac{f(r_1) - f(r_2)}{r_1 - r_2} = f'(r^*), \quad \frac{g(r_1) - g(r_2)}{r_1 - r_2} = g'(r^*), \quad \frac{h(r_1) - h(r_2)}{r_1 - r_2} = h'(r^*).$$

Step 3: Applying the Mean Value Theorem for Vector-Valued Functions The Mean Value Theorem (MVT) can be extended to vector-valued functions. Specifically, if $\mathbf{F}(r)$ is continuous on $[r_1, r_2]$ and differentiable on (r_1, r_2) , then there exists some $r^* \in (r_1, r_2)$ such that:

$$\mathbf{F}(r_2) - \mathbf{F}(r_1) = \mathbf{F}'(r^*)(r_2 - r_1).$$

Step 4: Expressing the Derivative and Rearranging the Equation Recall that the derivative $\mathbf{F}'(r)$ is given by:

$$\mathbf{F}'(r) = (f'(r), g'(r), h'(r)) = \nabla\mathbf{F}(r).$$

Rewriting the equation from Step 3:

$$\mathbf{F}(r_1) - \mathbf{F}(r_2) = (r_1 - r_2)\mathbf{F}'(r^*).$$

Step 5: Confirming the Equality Holds Component-Wise This equation holds for each component:

$$\begin{aligned} f(r_1) - f(r_2) &= (r_1 - r_2)f'(r^*), \\ g(r_1) - g(r_2) &= (r_1 - r_2)g'(r^*), \\ h(r_1) - h(r_2) &= (r_1 - r_2)h'(r^*). \end{aligned}$$

Figure 15: An example data instance for **DR.** category (2).

Step 6: Concluding the Proof Since all components satisfy the Mean Value Theorem at the same point r^* , we have:

$$\mathbf{F}(r_1) - \mathbf{F}(r_2) = (r_1 - r_2)\nabla\mathbf{F}(r^*),$$

where $r^* \in (r_1, r_2)$.

Step 7: Final Statement Therefore, for any r_1 and r_2 , there exists an r^* between them such that the original equation holds, which completes the proof.

3. Reason

MVT cannot be used in 3D space, which is a deception or trap, and the counterexample is a good way to show that the proposition may not hold. But we changed it to a correct proof by extending MVT to vector-valued functions and ignoring the restrictions of MVT

Figure 16: An example data instance for **DR.** category (3).

1. Process 1

Question: Three pencils and a jumbo eraser cost \$1.24. Five pencils and a jumbo eraser cost \$1.82. No prices include tax. In cents, what is the cost of a pencil?

Step 1: Let's call the price of a pencil p and the price of a jumbo eraser e . Then we can write two equations.

Step 2: We have $3p + e = 1.24$ and $5p + e = 1.82$.

Step 3: To solve this system, let's subtract the first equation from the second equation. This will eliminate e .

Step 4: $5p + e - 3p - e = 1.82 - 1.24$.

Step 5: This simplifies to $2p = 0.58$. So $p = 0.29$.

Step 6: That means a pencil costs 29 cents.

Answer

29

2. Process 2

Question: Three pencils and a jumbo eraser cost \$1.24. Five pencils and a jumbo eraser cost \$1.82. No prices include tax. In cents, what is the cost of a pencil?

Step 1: I need to solve this problem to find the cost of a pencil in cents. The problem gives me two scenarios involving pencils and jumbo erasers, along with their total costs. Let's first understand what's being asked.

Step 2: The problem states:

- Three pencils and one jumbo eraser cost \$1.24.

- Five pencils and one jumbo eraser cost \$1.82.

And it asks for the cost of a pencil in cents.

Step 3: Since the prices are given in dollars, I should convert them to cents to make calculations easier because the answer is required in cents.

- \$1.24 is equal to 124 cents.

- \$1.82 is equal to 182 cents.

So, I can rephrase the problem as:

- 3 pencils + 1 jumbo eraser = 124 cents

- 5 pencils + 1 jumbo eraser = 182 cents

Step 4: Now, I need to find the cost of one pencil. Let's denote:

- Let p be the cost of one pencil in cents.

- Let e be the cost of one jumbo eraser in cents.

So, I can write the equations based on the given information:

Equation 1: $3p + e = 124$

Equation 2: $5p + e = 182$

Step 5: I have a system of two equations with two variables (p and e). To find the value of p , I can use the method of elimination or substitution. I think elimination might be straightforward here.

Step 6: If I subtract Equation 1 from Equation 2, the e (eraser cost) will cancel out, and I can solve for p .

So, Equation 2 - Equation 1:

$(5p + e) - (3p + e) = 182 - 124$

Simplifying:

$5p + e - 3p - e = 58$

Which reduces to:

$2p = 58$

Step 7: Now, solving for p :

$p = 58 / 2$

$p = 29$ cents

Figure 17: An example data instance for **MS**. category (1).

Step 8: So, the cost of one pencil is 29 cents. But to make sure, I should check if this makes sense with both original equations.

Step 9: Plugging $p = 29$ into Equation 1:

$$3(29) + e = 124$$

$$87 + e = 124$$

$$e = 124 - 87$$

$$e = 37 \text{ cents}$$

Step 10: Now, plugging $p = 29$ and $e = 37$ into Equation 2 to verify:

$$5(29) + 37 = 145 + 37 = 182 \text{ cents, which matches the second equation.}$$

Step 11: Everything checks out, so the cost of a pencil is indeed 29 cents.

Final Answer

29

Figure 18: An example data instance for **MS.** category (2).

1. System Prompt

Task Description You are a helpful AI assistant who is very good at reasoning and data construction. Now I want to test the ability of process-level reward models to judge whether a step within the reasoning process is correct. To do this, please help me build flawed cases by introducing specific types of errors into a given reasoning process.

You will be provided with:

1. A mathematical problem.
2. Its standard correct answer.
3. A correct step-by-step reasoning process used to solve it.

Your task is to modify the question, adjust one or more steps, or introduce additional steps into the original process chain to create a reasoning process that appears plausible but is incorrect. The objective is to simulate flawed solutions by incorporating the specified error detailed after "Error Type to Introduce".

Error Type to Introduce

Redundancy refers to a process that is not the most concise or efficient, as it includes one or more redundant steps that can be removed without affecting the correctness of the overall solution path. For example, if $A \rightarrow B$ represents a correct inference chain, your task is to introduce one or more redundant steps $C = \{c | c \text{ is redundant}\}$ and reformulate the solution chain as $A \rightarrow C \rightarrow B$.

Formatting Instructions

After making the modifications, provide the following structured output:

```
{
  "original_question": "The original mathematical problem.",
  "modified_question": "The modified problem or original problem",
  "original_process": ["original_step 1", "original_step 2", ...],
  "modified_process": ["modified_step 1", "modified_step 2", ...],
  "modified_steps": [1, 5, 7, ...],
  "error_steps": [5, 6, ...],
  "reason": "Explanation for the changes."
}
```

Detailed Requirements:

1. original_question: A string representing the original mathematical problem as provided.
2. modified_question: A string representing the modified problem after your changes. If the problem remains the same, you can copy the original question.
3. original_process: A non-empty list of strings representing the original reasoning steps provided as input.
4. modified_process: A non-empty list of strings representing the reasoning process after your modifications.
5. modified_steps: A non-empty list of integers indicating the indexes of all modified steps. Indexing starts at 1.
6. error_steps: A non-empty list of integers representing the steps that contain hallucinations or errors. These should also be part of modified_steps.
7. reason: A clear explanation of the modifications made, why they were introduced, and how they align with the specified error types.

Notes:

1. Ensure all lists are non-empty.
2. Use the LaTeX format for all mathematical symbols (e.g., x^2 for x squared). Do not use Unicode symbols such as $\u2248$ or $\u00f7$.
3. Ensure the JSON object is well-formed, with proper escaping for special characters like backslash n (e.g., use backslash backslash n for newlines).
4. All indexes start from 1, that is, the first step's index is 1, not 0.
5. You can choose to modify the question or not, if the question remains the same, you can copy the original question. But if the question is modified, ensure that the steps are judged based on the modified question.
6. Please give the original process as provided by the prompt, do not modify it.

Figure 19: Prompts used for test case construction during data generation (1).

2. Few Shots

User 1

Question:

If $x^2 - x - 1 = 0$, what is the value of $x^3 - 2x + 1$?

Steps:

Step 1. Let's start with the first equation and see if we can solve for x.

Step 2. We can use the quadratic formula to find x.

Step 3. The quadratic formula is $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$.

Step 4. Substituting $a = 1$, $b = -1$, and $c = -1$, we get $x = \frac{1 \pm \sqrt{5}}{2}$.

Step 5. Since $x^2 - x - 1 = 0$, we replace x^2 with $x + 1$ in $x^3 - 2x + 1$.

Step 6. Substituting, $x^3 - 2x + 1 = x(x + 1) - 2x + 1$.

Step 7. Simplify: $x^3 - 2x + 1 = (x^2 + x) - 2x + 1$.

Step 8. Substitute $x^2 = x + 1$: $x^3 - 2x + 1 = (x + 1 + x) - 2x + 1$.

Step 9. Simplify: $x^3 - 2x + 1 = 2$.

Answer

2

Assistant 1

```
{
  "original_question": "If  $x^2 - x - 1 = 0$ , what is the value of  $x^3 - 2x + 1$ ?",
  "modified_question": "If  $x^2 - x - 1 = 0$ , what is the value of  $x^3 - 2x + 1$ ?",
  "original_process": [
    "Let's start with the first equation and see if we can solve for x.",
    "We can use the quadratic formula to find x.",
    "The quadratic formula is  $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$ ",
    "Substituting  $a = 1$ ,  $b = -1$ , and  $c = -1$ , we get  $x = \frac{1 \pm \sqrt{5}}{2}$ .",
    "Since  $x^2 - x - 1 = 0$ , we replace  $x^2$  with  $x + 1$  in  $x^3 - 2x + 1$ .",
    "Substituting,  $x^3 - 2x + 1 = x(x + 1) - 2x + 1$ .",
    "Simplify:  $x^3 - 2x + 1 = (x^2 + x) - 2x + 1$ .",
    "Substitute  $x^2 = x + 1$ :  $x^3 - 2x + 1 = (x + 1 + x) - 2x + 1$ .",
    "Simplify:  $x^3 - 2x + 1 = 2$ ."
  ],
  "modified_process": [
    "Let's start with the first equation and see if we can solve for x.",
    "We can use the quadratic formula to find x.",
    "The quadratic formula is  $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$ ",
    "Substituting  $a = 1$ ,  $b = -1$ , and  $c = -1$ , we get  $x = \frac{1 \pm \sqrt{5}}{2}$ .",
    "Since  $x^2 - x - 1 = 0$ , we replace  $x^2$  with  $x + 1$  in  $x^3 - 2x + 1$ .",
    "Substituting,  $x^3 - 2x + 1 = x(x + 1) - 2x + 1$ .",
    "Simplify:  $x^3 - 2x + 1 = (x^2 + x) - 2x + 1$ .",
    "Now let's verify that  $x^2 = x + 1$  by re-solving  $x^2 - x - 1 = 0$  again.",
    "Substitute  $x^2 = x + 1$ :  $x^3 - 2x + 1 = (x + 1 + x) - 2x + 1$ .",
    "Simplify:  $x^3 - 2x + 1 = 2$ ."
  ],
  "modified_steps": [8],
  "error_steps": [8],
  "reason": "Step 8 re-solves  $x^2 - x - 1 = 0$  unnecessarily, introducing redundancy. This step does not provide new information, as  $x^2 = x + 1$  was already established in The reasoning is still correct but less efficient."
}
```

User 2

Figure 20: Prompts used for test case construction during data generation (2).

Question:

In 1992, a scoop of gelato could be purchased in Italy for 1200 lire. The same gelato would have cost \$1.50 in the U.S. At the equivalent exchange rate between the lire and the dollar, how many dollars would be equivalent to 1,000,000 lire?

Steps :

Step 1. Let's call the number of dollars x. So, the problem is asking us to solve the equation $x = 1000000 * \frac{1.50}{1200}$.

Step 2. There are two ways to solve this equation. We could either divide 1000000 by 1200 and then multiply the result by 1.50 or we could divide 1.50 by 1200 and then multiply the result by 1000000.

Step 3. I think the second way will be easier.

$\frac{1.50}{1200} = 0.00125$.

Step 4. That's true. Now we just have to solve the equation $x = 1000000 * 0.00125$.

Step 5. So $x = 1250$.

Step 6. That's the final answer.

Step 7. Right. So 1,000,000 lire is equivalent to \$1250.

Answer

1250

Assistant 2

```
{
  "original_question": "In 1992, a scoop of gelato could be purchased in Italy for 1200 lire. The same gelato would have cost $1.50 in the U.S. At the equivalent exchange rate between the lire and the dollar, how many dollars would be equivalent to 1,000,000 lire?",
  "modified_question": "In 1992, a scoop of gelato could be purchased in Italy for 1200 lire. The same gelato would have cost $1.50 in the U.S. At the equivalent exchange rate between the lire and the dollar, how many dollars would be equivalent to 1,000,000 lire?",
  "original_process": [
    "Let's call the number of dollars x. So, the problem is asking us to solve the equation  $x = 1000000 * \frac{1.50}{1200}$ .",
    "There are two ways to solve this equation. We could either divide 1000000 by 1200 and then multiply the result by 1.50 or we could divide 1.50 by 1200 and then multiply the result by 1000000.",
    "I think the second way will be easier.  $\frac{1.50}{1200} = 0.00125$ .",
    "That's true. Now we just have to solve the equation  $x = 1000000 * 0.00125$ .",
    "So  $x = 1250$ .",
    "That's the final answer.",
    "Right. So 1,000,000 lire is equivalent to $1250."
  ],
  "modified_process": [
    "Let's call the number of dollars x. So, the problem is asking us to solve the equation  $x = 1000000 * \frac{1.50}{1200}$ .",
    "There are two ways to solve this equation. We could either divide 1000000 by 1200 and then multiply the result by 1.50 or we could divide 1.50 by 1200 and then multiply the result by 1000000.",
    "I think the second way will be easier.  $\frac{1.50}{1200} = 0.00125$ .",
    "That's true. But let's confirm the computation again:  $\frac{1.50}{1200} = 0.00125$ .",
    "Now we just have to solve the equation  $x = 1000000 * 0.00125$ .",
    "So  $x = 1250$ .",
    "Let's double-check by performing the multiplication again:  $1000000 * 0.00125 = 1250$ .",
    "That's the final answer.",
    "Right. So 1,000,000 lire is equivalent to $1250."
  ],
  "modified_steps": [4, 7, 8],
}
```

Figure 21: Prompts used for test case construction during data generation (3).

```
"error_steps": [4, 7, 8],  
"reason": "Step 4 unnecessarily repeats the computation of  $\frac{1.50}{1200}$ , which was already completed in Step 7 redundantly re-checks the multiplication, adding no new information. These redundant steps make the reasoning process less concise without affecting correctness."  
}
```

Figure 22: Prompts used for test case construction during data generation (4).

1. System Prompt

You are a mathematical reasoning evaluator. Your task is to analyze mathematical problem-solving steps and provide structured assessments in JSON format.

For each solution step, you need to evaluate two aspects:

1. Validity Score (-1 to +1):

- * +1: Completely correct mathematical reasoning
- * 0: Partially correct with some mistakes
- * -1: Completely incorrect
- * Use any value in between to indicate varying degrees of correctness

2. Redundancy Score (-1 to +1):

- * -1: Critical step, absolutely necessary for the solution
- * 0: Moderately important step
- * +1: Completely redundant, can be omitted
- * Use any value in between to indicate varying degrees of redundancy

Requirements:

- Evaluate each step independently
- Provide scores as floating-point numbers
- Return results in strict JSON format: "validity": [scores], "redundancy": [scores]
- Ensure both arrays have the same length as the number of steps
- Maintain mathematical rigor in your evaluation
- Consider mathematical accuracy, logical coherence, and solution efficiency

Example output format:

```
{ "validity": [0.8, -0.5, 1.0], "redundancy": [-1.0, 0.3, 0.7]}
```

2. Few Shots

User 1

Question:

In 1992, a scoop of gelato could be purchased in Italy for 1200 lire. The same gelato would have cost \$1.50 in the U.S. At the equivalent exchange rate between the lire and the dollar, how many dollars would be equivalent to 1,000,000 lire?

Solution:

Step 1. Let's call the number of dollars x . So, the problem is asking us to solve the equation $x = 1000000 * \frac{1.50}{1200}$.

Step 2. There are two ways to solve this equation. We could either divide 1000000 by 1200 and then multiply the result by 1.50 or we could divide 1.50 by 1200 and then multiply the result by 1000000.

Step 3. I think the second way will be easier. $\frac{1.50}{1200} = 0.00125$.

Step 4. That's true. Now we just have to solve the equation $x = 1000000 * 0.00125$.

Step 5. So $x = 1250$.

Step 6. So $x = 1250$.

Step 7. That's the final answer.

Step 8. Right. So 1,000,000 lire is equivalent to \$1250.

Answer

1250

Assistant 1

```
{\ "validity\": [1.0, 0.5, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0], \"redundancy\": [-1.0, 0.5, -0.5, -0.5, -1.0, 1.0, 0.5, 1.0]}
```

User 2

Figure 23: Prompts used for evaluating generative LLMs (1).

Question:

Four points, A , B , C , and D , are chosen randomly and independently on the circumference of a circle. What is the probability that segments AB and CD intersect?

Solution:

Step 1. Let's think about what it means for the segments to intersect.

Step 2. Whether they intersect is entirely a function of the order of B, C, D on the circle, when we look at it counterclockwise and consider A to come first.

Step 3. Exactly. So let's consider the different possible orders of B, C, D .

Step 4. The number of possible orders is $3! = 6$.

Step 5. Yes, that's right. Now, we need to figure out how many of these result in the segments intersecting.

Step 6. If C is in between B and D , then the segments will intersect.

Step 7. If B is in between C and D , then the segments will intersect.

Step 8. Right. That's two of the possible orders.

Step 9. So, the probability that the segments intersect is $\frac{2}{6} = \frac{1}{3}$.

Step 10. And that's our answer.

Answer

1/3

Assistant 2

```
{"validity": [1.0, 1.0, 1.0, 1.0, 1.0, -1.0, 0.8, 1.0, 1.0, 1.0], \"
  redundancy\": [-1.0, -0.5, -0.5, -0.5, -0.5, -0.5, 0.5, -0.5, -1.0, 1.0]}
```

Figure 24: Prompts used for evaluating generative LLMs (2).