

Dynamic and Generalizable Process Reward Modeling

Zhangyue Yin[◇] Qiushi Sun[♡] Zhiyuan Zeng[◇]
Qinyuan Cheng[◇] Xipeng Qiu^{◇†} Xuanjing Huang^{◇†}

[◇]Fudan University [♡]The University of Hong Kong
{yinzy21, cengzy23, chengqy21}@m.fudan.edu.cn
qiushisun@connect.hku.hk {xpqiu, xjhuang}@fudan.edu.cn

Abstract

Process Reward Models (PRMs) are crucial for guiding Large Language Models (LLMs) in complex scenarios by providing dense reward signals. However, existing PRMs primarily rely on heuristic approaches, which struggle with cross-domain generalization. While LLM-as-judge has been proposed to provide generalized rewards, current research has focused mainly on feedback results, overlooking the meaningful guidance embedded within the text. Additionally, static and coarse-grained evaluation criteria struggle to adapt to complex process supervision. To tackle these challenges, we propose Dynamic and Generalizable Process Reward Modeling (DG-PRM), which features a reward tree to capture and store fine-grained, multi-dimensional reward criteria. DG-PRM dynamically selects reward signals for step-wise reward scoring. To handle multifaceted reward signals, we pioneeringly adopt Pareto dominance estimation to identify discriminative positive and negative pairs. Experimental results show that DG-PRM achieves stunning performance on prevailing benchmarks, significantly boosting model performance across tasks with dense rewards. Further analysis reveals that DG-PRM adapts well to out-of-distribution scenarios, demonstrating exceptional generalizability.

“Judgements prevent us from seeing the good that lies beyond appearances.”

–Wayne Dyer

1 Introduction

The Process Reward Model (PRM) plays a crucial role in enhancing the performance of Large Language Models (LLMs) (Uesato et al., 2023; Yuan et al., 2024; Zhang et al., 2025). Unlike outcome reward models (ORM), which provide feedback solely on the final objective, PRM focuses on rewarding the intermediate steps or processes involved in task execution (Setlur et al., 2025). This

Answer:	Judgement:
Step 1 The gravitational force provides the centripetal force for circular motion: $F_g = \frac{GMm}{r^2} = \frac{mv^2}{r}$, where M is Earth's mass, m is the satellite's mass, v is orbital speed, and r is orbital radius.	Step 1 (✓) The step 1 correctly explains the gravitational force providing the centripetal force and establishes the relationship between orbital speed, orbital radius, and the mass of the Earth. [[Correct]]
Step 2 Rearranging the equation gives $v = \sqrt{\frac{GM}{r}}$, showing orbital speed decreases as r increases.	Step 2 (✓) The conclusion that the orbital speed decreases as the orbital radius increases is accurate. [[Correct]]
Step 3 If v decreases, the gravitational force exceeds the required centripetal force, causing the satellite to spiral inward to a smaller orbit.	Step 3 (✗) The error occurs in the third step. The misconception here is that the satellite would spiral inward when its speed decreases. [[Incorrect]]
Step 4 Using Kepler's third law, $T^2 \propto r^3$, a larger r implies a longer orbital period T .	Step 4 (✓) The step correctly applies Kepler's third law and concludes that a larger orbital radius results in a longer orbital period. [[Correct]]
Step 5 Applying the escape velocity formula gives $v = \sqrt{\frac{2GM}{r}}$, the new speed corresponds to $r_{new} = \frac{2GM}{v^2}$.	Step 5 (✗) The mistake happened in the fifth step. The escape velocity formula $v = \sqrt{\frac{2GM}{r}}$ was incorrectly applied instead of the correct orbital velocity formula $v = \sqrt{\frac{GM}{r}}$. [[Incorrect]]
Step 6 The total mechanical energy $E = -\frac{GMm}{2r}$ increases as v decreases, so r decreases to conserve energy.	Step 6 (✗) The claim that the orbital radius decreases to conserve energy is incorrect. In a stable orbit, as the satellite's speed decreases, the orbital radius must increase. [[Incorrect]]

Figure 1: An example from MT-Bench, where LLM judgments contain rich details, as highlighted in bold, but prevailing PRM methods fail to fully exploit.

approach provides more dense reward signals (Li and Li, 2024), especially in complex problems where the model needs to reason, analyze, and explore different solution strategies (Wei et al., 2022). By evaluating intermediate steps, PRM is essential for improving the model's ability to tackle intricate tasks (Lightman et al., 2024).

PRM can be broadly classified into heuristic and generative categories. Heuristic process rewards (Wang et al., 2024a,c) rely on manually crafted criteria to assess the relevance of intermediate steps to the final answer. While heuristic rewards have significantly advanced the model's reasoning capabilities, they suffer from several limitations (Zhang et al., 2025). Specifically, they often require objective, reference-based answers (Luo et al., 2024), which are difficult to obtain in complex scenarios where evaluation criteria fluctuate. Furthermore, PRMs trained using heuristic rewards often exhibit poor generalization and even be susceptible to reward hacking (Weng, 2024; Wen et al., 2025). On the other hand, generative process rewards utilize LLMs to replace human annotation by labeling each intermediate step as correct or

incorrect, offering positive or negative feedback accordingly (Mahan et al., 2024; Cao et al., 2024). While generative rewards capitalize on the LLM’s ability to generate responses, the evaluation still relies on a fixed set of standards (Ling et al., 2023), limiting its adaptability across diverse domains.

Furthermore, prevailing LLM-as-Judge methods only utilize final feedbacks (*e.g.*, correct/incorrect) and overlook valuable information encapsulated in the process (Kwon et al., 2023; Gao et al., 2023b), such as error severity and the type of mistakes. As illustrated in Figure 1, we observe that LLM feedback contains rich details and guidance information. However, the prevalent approach assigns a uniform negative reward for incorrect labels, neglecting the diversity and severity of errors.

In this paper, we identify two key limitations in current process of reward construction: (1) the use of fixed evaluation criteria, and (2) the reliance on uniform negative rewards, which fail to capture the diversity and severity of errors, limiting adaptability and generalizability in process reward. To address these challenges, we introduce Dynamic and Generalizable Process Reward Modeling (DG-PRM), a novel framework designed to automatically construct and precisely allocate process rewards. We propose the use of a reward tree to store multifaceted evaluation criteria extracted from LLM judgments. It selects the most step-wise relevant criteria during evaluation, thus making DG-PRM excel in cross-domain generalization. We also introduce Pareto dominance estimation to select positive and negative pairs from a diverse set of reward signals, providing clear optimization objectives. Experimental results demonstrate that DG-PRM achieves state-of-the-art performance on PRMBENCH, showcasing superior PRM capabilities. By offering contextually appropriate reward signals, DG-PRM significantly improves LLM performance across a wide range of tasks. Furthermore, compared to the LLM-as-Judge approach, DG-PRM demonstrates enhanced training efficiency and better generalization to out-of-distribution scenarios.

Our main contributions are listed below:

- We introduce DG-PRM, an automated framework designed to construct dynamic and generalizable process rewards, optimizing the utility of LLM feedback.
- To handle diverse and complex rewards, we introduce a novel reward tree to dynamically

capture and leverage appropriate criteria to each evaluation step.

- We propose the use of Pareto dominance estimation to identify positive and negative pairs from multifaceted reward signals, thereby providing clearer optimization objectives.
- DG-PRM significantly boosts LLM performance across a wide range of tasks by offering precise, fine-grained process rewards, while demonstrating high training efficiency and exceptional generalizability.

2 Related Work

Outcome Reward Model. Reward models are designed to capture human preferences and automate the evaluation of model outputs (Ouyang et al., 2022; Kaufmann et al., 2024; Sun et al., 2024a). Outcome Reward Model (ORM) has been applied across a broad range of domains, including safety (Dai et al., 2024), mathematical problem-solving (Cobbe et al., 2021; Yang et al., 2024b), and code generation (Dou et al., 2024; Sun et al., 2024b). Recent studies, such as Wang et al. (2024b), have further expanded the reward signal to encompass diverse dimensions, including helpfulness, correctness, coherence, complexity, and verbosity. ORM can be categorized into discriminative models (Stiennon et al., 2020; Ouyang et al., 2022) and generative models (Mahan et al., 2024). Discriminative reward models typically add a classification head to assess the quality of inputs (Gao et al., 2023a; Chen et al., 2024a), whereas generative reward models leverage the language generation capabilities of LLMs to evaluate outputs (Zhu et al., 2024; Li et al., 2024c). Zheng et al. (2023) demonstrate that LLMs can provide scalable and explainable rewards that exhibit high alignment with human preferences.

Process Reward Model. As LLMs are increasingly required to handle complex tasks, they often need more tokens to reason effectively (Wei et al., 2022; Snell et al., 2024). As the length of LLM outputs increases, ORM struggles to fully evaluate the coherence and correctness of the output (Luo et al., 2024). Thus, building PRMs with dense reward signals emerges as a solution (Lightman et al., 2024). Current approaches of building PRMs tend to focus on objective domains, such as mathematics, which have clear, definitive answers (Guan et al., 2025), aiming to improve LLM performance in mathematical problem-solving (Uesato et al., 2022; Yuan

et al., 2024; Cui et al., 2025). Wang et al. (2024a) model the correctness of each step’s output as a process reward, while Wang et al. (2024c) adopt a softer approach by incorporating the likelihood of the current step’s output being correct as a process reward. However, these heuristic methods limit the generalizability, particularly in cases where there is no clear, unique answer. For example, in scientific tasks, diverse and complex reward signals should be considered, and various output components may require attention to distinct reward signals (Wu et al., 2024). Therefore, constructing dynamic and diverse process rewards is a desideratum.

Reward Signal. The design of reward signals is a crucial component in RL (Sutton, 2018). Reward signals can generally be classified into three types: human-annotated (Bai et al., 2022a), rule-based (Glaese et al., 2022), and AI-feedback (Lee et al., 2024b). Human-annotated reward signals require expert labeling and verification, which can be costly and time-consuming (Lightman et al., 2024). Mu et al. (2024) argue that human annotations often fail to accurately convey the intended behaviors to annotators, which complicates the conversion of desired outcomes into specific rules. Rule-based systems, such as parsing, use compilers to perform syntactic analysis for translating source code into executable binary code. However, such approaches are domain-specific and cannot easily generalize to other areas. With the increasing capabilities of LLMs, research has shifted toward utilizing AI feedback for reward generation (Bai et al., 2022b, 2023; Li et al., 2024b). Kwon et al. (2023) observe that LLMs, when used as a proxy reward function, significantly improve the alignment of rewards with user objectives. In the context of code, McAleese et al. (2024) show that AI models help identify more bugs than human contractors. Cao et al. (2024) demonstrate that dense reward signals provided by LLMs effectively improve the performance of policy models. However, current research (Gao et al., 2024; Chen et al., 2024b; Ling et al., 2023) typically focuses on scoring or ranking outputs, which neglects the rich guiding information present in LLM-generated texts.

3 Preliminary

Given the complexity of process rewards, we design a reward tree to store high-quality reward criteria at different levels of granularity. The reward tree dynamically selects the most appropriate criteria

at each step, and the corresponding reward scores are computed based on these criteria. To tackle the optimization challenges posed by the diversity of reward signals, we introduce Pareto dominance estimation and select Pareto-optimal samples. The detailed process is illustrated in Figure 2.

3.1 Automatic Process Reward Design

Given an input x , the model output \hat{y} is evaluated against a reference y or a set of evaluation criteria \mathcal{C} . The output \hat{y} is categorized into two sets, \hat{y}_+ and \hat{y}_- , such that:

$$\hat{y}_+ = y \quad \text{and} \quad \hat{y}_- \neq y \quad \text{or} \quad \mathcal{C}(\hat{y}_+) > \mathcal{C}(\hat{y}_-) \quad (1)$$

Function \mathcal{J} then analyzes why \hat{y}_+ is superior to \hat{y}_- , yielding a set of reward criteria:

$$\mathcal{J}(x, \hat{y}_+, \hat{y}_-) = \{r_1, r_2, r_3, \dots\} \quad (2)$$

After applying this process to all samples in the dataset, an automated validator filters out low-quality and unreasonable reward criteria (detailed in Appendix B.2). This results in a set of m reward criteria: $\{r_1, r_2, r_3, \dots, r_m\}$. Each reward criterion r_i is then mapped to a d -dimensional vector space using the embedding function \mathcal{V} :

$$\mathcal{V}(r_i) = [v_i^{(1)}, v_i^{(2)}, \dots, v_i^{(d)}] \quad (3)$$

Next, we construct a reward tree \mathcal{T} using hierarchical clustering with the distance function $\mathcal{D}(v_i, v_j)$, which is defined as:

$$\mathcal{D}(v_i, v_j) = 1 - \frac{v_i \cdot v_j}{\|v_i\| \|v_j\|} \quad (4)$$

The reward tree is generated as follows:

$$\mathcal{T} = \mathcal{H}(\{v_1, v_2, \dots\}, \mathcal{D}) \quad (5)$$

To reduce redundancy among tree nodes, we define a threshold ξ such that when $\mathcal{D}(v_i, v_j) \leq \xi$, we merge r_i and r_j . To facilitate dynamic updates to the reward tree, we adopt an incremental hierarchical clustering algorithm, minimizing the cost of introducing new rewards. Based on the constructed hierarchical relationship, we annotate the upper-level coarse-grained reward criteria as r^{parent} and the lower-level fine-grained reward criteria as r^{child} . The final reward tree \mathcal{T} is then expressed as:

$$\mathcal{T} = (\{r_1^{\text{parent}}, r_2^{\text{parent}}, \dots\}, \{r_1^{\text{child}}, \dots, r_m^{\text{child}}\}) \quad (6)$$

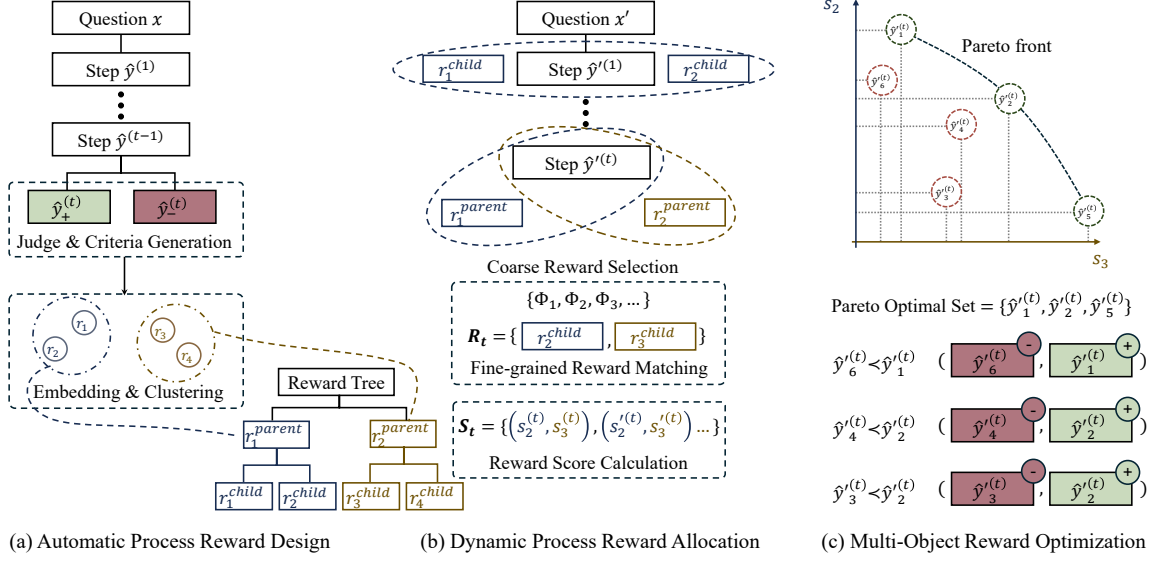


Figure 2: Overview of DG-PRM. DG-PRM consists of three main steps: (a) Automatic Process Reward Design, which constructs reward criteria using positive and negative sample pairs \hat{y}_+ and \hat{y}_- , maps these criteria into a feature space, and builds a reward tree via hierarchical clustering; (b) Dynamic Process Reward Allocation, which dynamically selects both coarse-grained rewards r^{parent} and fine-grained rewards r^{child} from the reward tree in each step, computing the reward score based on the corresponding criteria; and (c) Multi-Objective Reward Optimization, which selects the Pareto-optimal $\hat{y}^{(t)}$ as the optimization target based on the computed reward scores.

3.2 Dynamic Process Reward Allocation

The dynamic process reward allocation uses the reward tree \mathcal{T} to provide step-wise reward signals. Given input data x and model output \hat{y} , the output is divided into n steps:

$$\hat{y} = \{\hat{y}^{(1)}, \hat{y}^{(2)}, \dots, \hat{y}^{(n)}\} \quad (7)$$

For each step t , the reward allocation dynamically assigns corresponding rewards. Specifically, for the t -th step $\hat{y}^{(t)}$, the function \mathcal{R} selects coarse-grained reward r^{parent} from the upper levels of reward tree \mathcal{T} , and the analysis function Φ further determines whether fine-grained reward criteria r^{child} are needed for evaluation based on $\hat{y}^{(t)}$ and the selected r^{parent} .

$$\{r_i^{\text{parent}}, \Phi(\hat{y}^{(t)}, r_i^{\text{parent}}), \dots\} = \mathcal{R}(\hat{y}^{(t)}, \mathcal{T}) \quad (8)$$

If fine-grained reward criteria r^{child} are needed, to ensure matching accuracy, we extract each evaluation criterion generated by Φ , denoted as $\phi_i = \Phi(\hat{y}^{(t)}, r^{\text{parent}})$. Similarly, we use the embedding function \mathcal{V} to map ϕ_i to the d -dimensional vector space, and compute the cosine distance $\delta_k^{(t)}$ between ϕ_i and each child node $r_k^{(t)}$ under r^{parent} .

$$\delta_k^{(t)} = \mathcal{D}(\mathcal{V}(\phi_i), \mathcal{V}(r_k^{(t)})) \quad (9)$$

For those $r_k^{(t)}$ satisfying $\delta_k^{(t)} \leq \zeta$, where ζ is a pre-defined threshold. Thus, the reward set for the t -th

step is:

$$\mathbf{R}_t = \{r_k^{(t)} | \delta_k^{(t)} \leq \zeta\} \quad (10)$$

We score each reward using a scoring function \mathcal{S} :

$$s_k^{(t)} = \mathcal{S}(\hat{y}^{(t)}, r_k^{(t)}), \quad r_k^{(t)} \in \mathbf{R}_t \quad (11)$$

For ensuring the non-independence of process rewards, the evaluation of step t incorporates the information from the preceding t steps, i.e., $\hat{y}^{(<t)}, s_k^{(<t)}$. To avoid excessively long contexts, only the last μ steps are retained, as detailed in Algorithm 1.

3.3 Multi-Object Reward Optimization

After assigning rewards and computing scores for the current step, we obtain a reward set \mathbf{R}_t . Given the diversity within \mathbf{R}_t , we select an optimal candidate from the score set \mathbf{S}_t as our optimization target. For any candidate $\hat{y}_i^{(t)}$, if there exists another candidate $\hat{y}_j^{(t)}$ such that for every reward $r_k^{(t)} \in \mathbf{R}_t$ the following holds:

$$s_k^{(t,i)} \geq s_k^{(t,j)}, \quad (12)$$

and there exists at least one index l for which

$$s_l^{(t,i)} > s_l^{(t,j)}, \quad (13)$$

then we say that $\hat{y}_i^{(t)}$ Pareto-dominates $\hat{y}_j^{(t)}$, denoted as $\hat{y}_i^{(t)} \succ \hat{y}_j^{(t)}$.

Algorithm 1 Dynamic Process Reward Allocation

Require: Step- t $\hat{y}^{(t)}$, Reward tree \mathcal{T} , Distance threshold ζ , Window size μ , Reward function \mathcal{R} , Analysis function Φ , Embedding function \mathcal{V} , Distance function \mathcal{D} , Score function \mathcal{S}

Ensure: Reward set \mathbf{R} , Score set \mathbf{S}

```
1: Initialize  $\mathbf{R} = \emptyset, \mathbf{S} = \emptyset$ 
2: for each timestep  $t = 1$  to  $n$  do
3:   if  $t - \mu \leq 0$  then
4:     Select all available previous steps up
     to  $t - 1$ , i.e.,  $\{\hat{y}^{(1)}, \hat{y}^{(2)}, \dots, \hat{y}^{(t-1)}\}$ 
5:   else
6:     Select the previous  $\mu$  steps, i.e.,
      $\{\hat{y}^{(t-\mu)}, \hat{y}^{(t-\mu+1)}, \dots, \hat{y}^{(t-1)}\}$ 
7:   end if
8:   Retrieve corresponding rewards and
     scores:  $\mathcal{I}_t = \{(\hat{y}^{(t-\mu)}, r_k^{(t-\mu)}, s_k^{(t-\mu)}), \dots\}$ 
9:   Add  $\mathcal{I}_t$  as supplementary information to
     reward function  $\mathcal{R}$ 
10:  Cal  $\phi_i = \Phi(\hat{y}^{(t)}, r_i^{\text{parent}})$ 
11:  for each  $\phi_i \in \mathcal{R}(\hat{y}^{(t)}, \mathcal{T}, \mathcal{I}_t)$  do
12:    for each child  $r_k^{(t)}$  of  $r_i^{\text{parent}}$  do
13:       $\delta_k^{(t)} = \mathcal{D}(\mathcal{V}(\phi_i), \mathcal{V}(r_k^{(t)}))$ 
14:      if  $\delta_k^{(t)} \leq \zeta$  then
15:        Add  $r_k^{(t)}$  to  $\mathbf{R}$ 
16:        Cal  $s_k^{(t)} = \mathcal{S}(r_k^{(t)}, \hat{y}^{(t)}, \mathcal{I}_t)$ 
17:        Add  $s_k^{(t)}$  to  $\mathbf{S}$ 
18:      end if
19:    end for
20:  end for
21: end for
```

Ideally, we aim to obtain an absolute optimal solution $\hat{y}^{(t)*}$ that satisfies

$$\forall j, k \quad s_k^{(t,*)} \geq s_k^{(t,j)} \quad (14)$$

When such an absolute optimum does not exist, we instead search for a non-dominated solution $\hat{y}_i^{(t)}$ for which no other candidate $\hat{y}_j^{(t)}$ exists that Pareto-dominates it. The set of all non-dominated solutions is referred to as the Pareto optimal set, denoted by \mathbf{U} . The scores corresponding to these solutions form the Pareto front. We then define the optimization target set as

$$\mathbf{V} = \{(\hat{y}_+^{(t)}, \hat{y}_-^{(t)}) \mid \forall \hat{y}_+^{(t)} \in \mathbf{U}, \exists \hat{y}_-^{(t)}, \hat{y}_+^{(t)} \succ \hat{y}_-^{(t)}\}. \quad (15)$$

We employ the DPO (Rafailov et al., 2023) optimization algorithm. Therefore, the optimization

objective for the policy model is as follows:

$$\mathcal{L}(\theta) = -\mathbb{E}_{x, \hat{y}^{(<t)}, (\hat{y}_+^{(t)}, \hat{y}_-^{(t)}) \in \mathbf{V}} [\log \sigma(\mathcal{F}_+ - \mathcal{F}_-)] \quad (16)$$

where,

$$\begin{aligned} \mathcal{F}_+ &= \beta \log \frac{\pi_\theta(\hat{y}_+^{(t)}|x; \hat{y}^{(<t)})}{\pi_{\text{ref}}(\hat{y}_+^{(t)}|x; \hat{y}^{(<t)})} \\ \mathcal{F}_- &= \beta \log \frac{\pi_\theta(\hat{y}_-^{(t)}|x; \hat{y}^{(<t)})}{\pi_{\text{ref}}(\hat{y}_-^{(t)}|x; \hat{y}^{(<t)})} \end{aligned} \quad (17)$$

Here, σ denotes the sigmoid function, π_θ represents the policy model to be optimized, and π_{ref} is the fixed reference model during training. The hyperparameter β controls the distance between the policy model and the reference model.

4 Experiments

4.1 Evaluation Datasets

To evaluate the efficiency of DG-PRM in process reward modeling, we use the PRMBENCH (Song et al., 2025) dataset. This benchmark, built on the PRM800K corpus, consists of 6k math problems across 9 distinct error categories, enabling comprehensive evaluation of process-level reward models.

To further assess DG-PRM’s effectiveness across a broad range of tasks, we incorporate three additional task sets representing different domains: general tasks, scientific tasks, and commonsense reasoning. These datasets include:

- *General*: MT-Bench (Zheng et al., 2023), Arena-Hard (Tianle et al., 2024)
- *Science*: QASC (Khot et al., 2020), ChemistryQA (Wei et al., 2021)
- *Commonsense*: StrategyQA (Geva et al., 2021), ARC-c (Clark et al., 2018)

A detailed description of the datasets, including answer types, training and test set distributions, and licensing information, can be found in Appendix A.

4.2 Experiment Settings

Implementation Details. For PRMBENCH, we use the same setup as in the *Prompted as Critic Models* configuration, which includes four open-source models: QwQ-Preview-32B (Team, 2024), DeepSeek-R1-Distill-Qwen-32B, DeepSeek-R1-Distill-Qwen-7B, and DeepSeek-R1 (Guo et al., 2025), along with two proprietary models: o1-mini (OpenAI, 2024c) and GPT-4o (OpenAI, 2024a). In *General*, *Science*, and *Commonsense*

Model Name	Overall	Simplicity			Soundness					Sensitivity			
		NR.	NCL.	Avg.	ES	SC.	DC.	CI	Avg.	PS	DR.	MS.	Avg.
<i>Open-source Discriminative Process Reward Model</i>													
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
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
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
<i>Prompted as Critic Models</i>													
o1-mini ^{†*}	68.8	65.6	63.7	64.6	74.5	67.7	73.8	72.3	72.1	61.8	64.8	100.0	75.5
GPT-4o [†]	66.8	57.0	62.4	59.7	72.0	69.7	70.7	71.1	70.9	62.5	65.7	99.2	75.8
QwQ-Preview-32B [†]	63.6	57.2	55.6	56.4	67.4	72.3	66.2	66.9	68.2	57.8	62.7	100.0	73.5
R1-Distill-Qwen-32B [†]	60.2	57.2	51.9	54.5	66.1	68.4	69.3	64.8	67.2	53.3	54.6	99.9	69.3
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	69.5	66.0	65.2	65.6	74.8	70.1	72.2	72.9	72.5	63.2	66.2	100.0	76.5
<i>Dynamic and Generalizable Process Reward Modeling</i>													
o1-mini	73.5	71.2	69.1	70.2	77.5	74.8	76.3	75.6	76.1	67.3	70.4	100.0	79.2
GPT-4o	72.3	66.1	69.0	67.6	75.9	73.2	74.7	76.4	75.1	66.8	70.9	100.0	79.2
QwQ-Preview-32B	70.0	63.2	65.4	64.3	72.4	74.3	72.9	74.5	73.5	64.5	67.9	100.0	77.5
R1-Distill-Qwen-32B	69.0	62.0	64.7	63.4	71.1	72.6	71.3	73.8	72.2	63.6	66.8	100.0	76.8
R1-Distill-Qwen-7B	65.2	60.4	62.1	61.3	69.8	68.1	69.5	72.1	69.9	62.4	64.2	100.0	75.5
DeepSeek-R1	76.5	74.1	72.3	73.2	80.1	77.5	78.9	79.4	79.0	71.0	74.3	100.0	81.8

Table 1: Performance comparison of DG-PRM and other strong baselines on PRMBENCH (Song et al., 2025) (PRM-Score %). The best results are highlighted in **bold**. [†] indicates results from the official leaderboard, and * denotes evaluation on a subset of 394 samples. The evaluation includes mainstream open-source and closed-source models for a fair comparison. Details of each dataset category and evaluation objectives are provided in Appendix A.

scenarios, we also include the models DeepSeek-R1-Distill-Qwen-1.5B and DeepSeek-R1-Distill-Qwen-14B. To construct the reward tree, we use the BAAI/bge-en-icl model (Li et al., 2024a) to build \mathcal{V} , with the dimensionality d set to 4096. For hierarchical clustering, we apply the BIRCH algorithm (Zhang et al., 1997) to create \mathcal{H} . Unless otherwise stated, we use DeepSeek-R1-Distill-Qwen-7B as the backbone model and set the hyperparameters as follows: the merge hyperparameter $\xi = 0.25$, the distance hyperparameter $\zeta = 0.2$, and step hyperparameter $\mu = 20$. For more detailed experimental procedures and hyperparameter analysis, refer to Appendix B.1 and Appendix C.

Baselines. For PRMBENCH, we use the official results as baselines. To further validate the effectiveness of DG-PRM, we compare it against the following setups in the *General*, *Science*, and *Commonsense* scenarios: (1) **Original**: The original model without any optimization. (2) **ORM**: We train a reward model to provide positive and negative feedback for complete outputs. (3) **Critic Models**: Building on Song et al. (2025), we guide LLMs with prompts to critique the solution in a step-by-

step manner. Following Mahan et al. (2024), we distinguish between **Critic Direct**, where models directly generate results, and **Critic CoT**, where the model first performs reasoning and analysis before presenting the answer. (4) **Human Annotation**: Training on human-annotated process labels as an upper bound.

We also include **Implicit PRM** (Yuan et al., 2024), a token-level reward baseline in our analysis, as detailed in Appendix D.2. We observe that a substantial amount of preference data is essential for training in order to effectively model DPO-equivalent rewards (Rafailov et al., 2024).

4.3 Main Results

PRMBENCH. Table 1 presents the results on PRMBENCH, where generative models demonstrate a clear advantage over discriminative models, a finding consistent with Zheng et al. (2024a). DG-PRM further significantly improves the process reward modeling capabilities of LLMs. In more challenging PS tasks, where PRMs must be capable of identifying missing conditions and prerequisite mistakes, R1-Distill-Qwen-7B shows a remarkable improvement from 45.6% to 62.4%. This notable

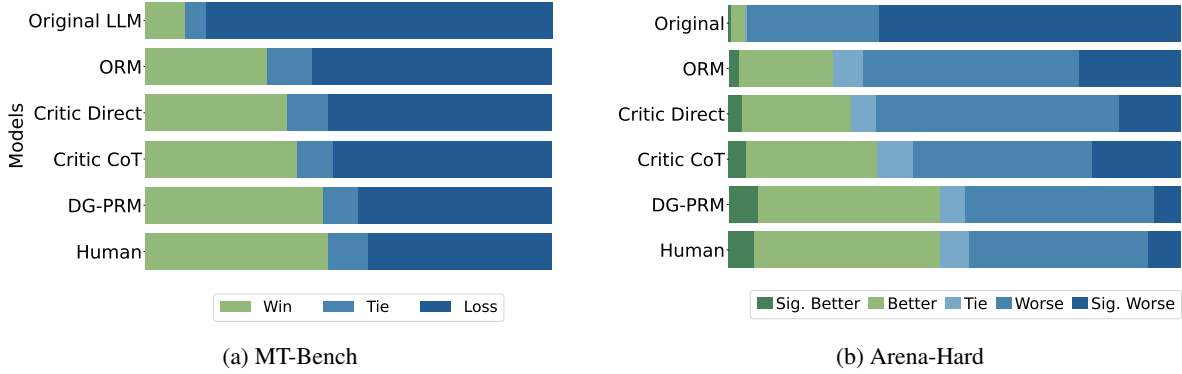


Figure 3: Performance comparison on (a) MT-Bench and (b) Arena-Hard. The R1-Distill-Qwen-7B model is used as the backbone, with GPT-4o (OpenAI, 2024a) serving as the judge model.

enhancement can be attributed to DG-PRM’s step-wise discrimination and dynamic reward allocation, which effectively increase error identification accuracy. Moreover, R1 demonstrates strong performance, which DG-PRM further enhances, improving overall accuracy from 69.5% to 76.5% as a new state-of-the-art. Notably, in sensitivity testing, the average score exceeds 80%. This can be attributed to R1’s complex reasoning generating higher-quality reward criteria, which DG-PRM effectively leverages to boost performance. We will further analyze this phenomenon in Appendix D.1.

General. Table 3 presents the results on MT-Bench and Arean-Hard. We observe that preference optimization significantly enhances the model’s win rate. The Critical Model outperforms ORM through step-wise analysis, achieving better performance. In Arean-Hard, using Critical CoT with process analysis demonstrates a clear advantage over Critical Direct, highlighting the importance of step-wise analysis during evaluation. DG-PRM notably improves the model’s win rate compared to the baseline, reaching levels close to Human Annotation, and even exceeding Human Annotation in the “Significantly Better” category on Arean-Hard. We attribute this to DG-PRM’s ability to integrate multi-dimensional reward information, resulting in more accurate process rewards.

Science and Commonsense. In Figure 4, we compare the performance of different methods across models with 1.5B, 7B, 14B, and 32B parameters. We observe that DG-PRM significantly enhances the performance of models at all parameter scales, demonstrating a clear advantage over other methods. Furthermore, DG-PRM achieves performance close to human annotation, even surpassing human-level performance on the Strate-

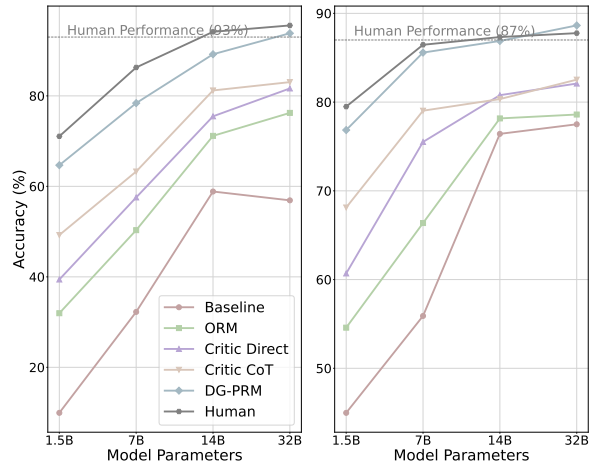


Figure 4: Performance comparison on (a) QASC and (b) StrategyQA. The evaluation includes models of different parameter scales: 1.5B, 7B, 14B, and 32B.

gyQA. Notably, the R1-Distill-Qwen-32B model trained with DG-PRM outperforms human performance on both the QASC and StrategyQA datasets. In StrategyQA, we observe a performance drop for the Critical CoT method, which even falls below the Critical Direct method in the 14B model. This is due to the implicit reasoning challenges in StrategyQA, where the model often struggles to determine the appropriate direction for analysis, particularly when the reasoning process is highly complex. In such cases, unoriented analysis can lead to erroneous judgments. DG-PRM, by providing relevant and explicit reward objectives at each step, effectively improves the judgments.

Generalization. In Figure 5, we analyze the generalization capabilities of different methods. We observe a significant performance drop for ORM and Critic Direct in out-of-domain (OOD) scenarios, with some results falling below baseline performance. To further analyze the Critical methods, we introduce Critic Instruction, which provides no

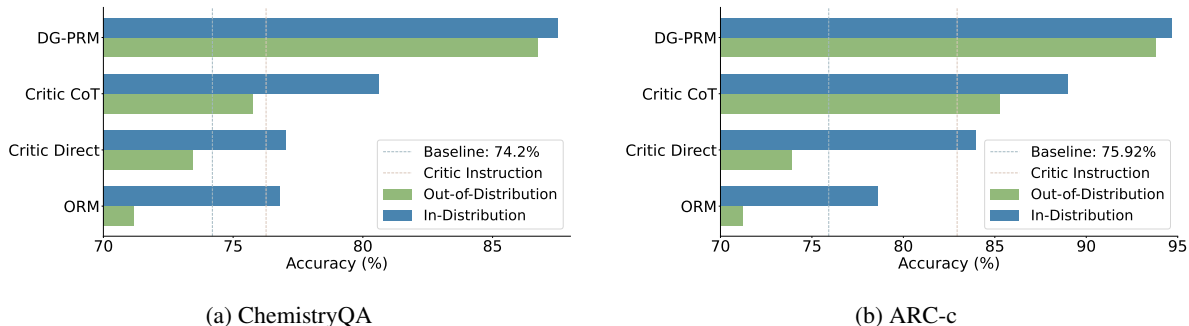


Figure 5: Generalization analysis on (a) ChemistryQA and (b) ARC-c. In the Out-of-Distribution setting, we provide process feedback for ChemistryQA samples using the process reward model constructed on QASC, and for ARC-c samples using the process reward model constructed on StrategyQA. **Critic-Instruction** refers to the approach where only instructions are used without providing any domain-specific exemplars.

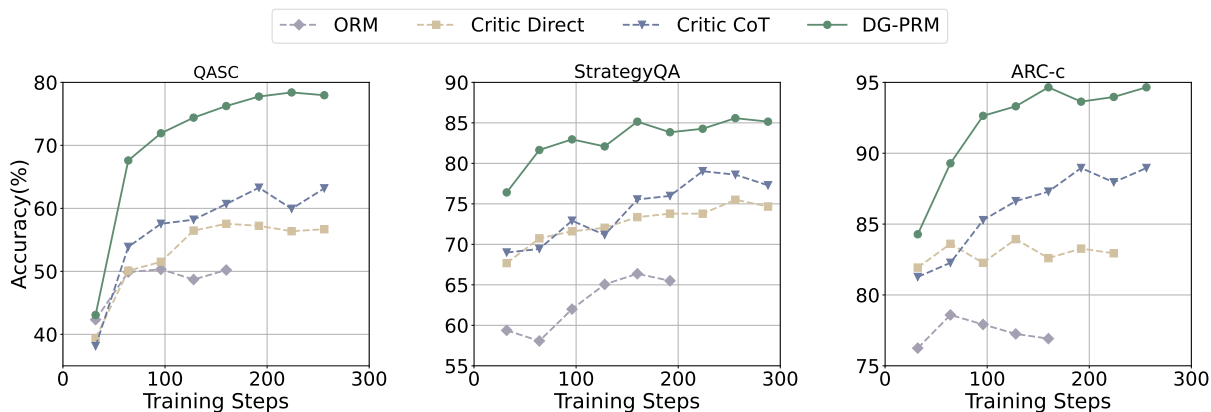


Figure 6: Accuracy(%) variation with training steps on the QASC, StrategyQA, and ARC-c datasets.

examples and only simple instructions to guide the LLM in generating process rewards. In the OOD setting, the Critical CoT method exhibits minimal performance degradation on the ARC-c, but a significant drop on ChemistryQA, even performing worse than Critic Instruction. This is attributed to the diverse problem distribution in QASC, which leads to confusion when transferring to the Chemistry domain. In contrast, DG-PRM, by selecting relevant and effective rewards, demonstrates outstanding domain generalization abilities. In Appendix D.5, we leverage the scalability of DG-PRM to construct a general reward tree for further analysis of its cross-domain generalization.

Training Efficiency. In Figure 6, we compare the accuracy of different methods as a function of training steps. We observe that DG-PRM demonstrates exceptional training efficiency across all datasets, achieving performance equivalent to Critic CoT with only 30% of the training steps. This efficiency is attributed to DG-PRM’s use of Pareto advantage estimation to select the most representative positive and negative samples, making it easier for the model to capture the differences between samples

and optimize the fitting objective. As a result, DG-PRM continuously improves model performance, significantly outperforming baseline methods.

5 Conclusion

In this paper, we first identify that current PRMs are often domain-specific and optimized for specific objectives, limiting their generalizability. The core challenge lies in the need for detailed, step-by-step supervision, making it difficult to standardize reward signals. To address this, we propose Dynamic and Generalizable Process Reward Modeling (DG-PRM), which utilizes a reward tree to store fine-grained reward signals and dynamically provides step-specific rewards. We also introduce Pareto dominance estimation to select positive and negative feedback from diverse reward signals. Experiments across benchmarks demonstrate that DG-PRM can deliver more accurate process reward signals. Notably, it achieves new state-of-the-art on PRMBENCH, and shows effectiveness in improving model performance across various domains, showcasing its exceptional generalizability by providing effective and precise reward signals.

Limitations

Expansion to Diverse Domains. Although we have demonstrated the effectiveness of DG-PRM across several datasets, including PRM-BENCH (Song et al., 2025), MT-Bench (Zheng et al., 2023), and Arena-Hard (Tianle et al., 2024), there is still significant potential for DG-PRM to expand into more domains. As LLMs are increasingly applied to various scientific fields, DG-PRM can be adapted to design domain-specific reward trees for fields such as drug discovery (Zheng et al., 2024b), disease diagnosis (Zhou et al., 2024), and weather forecasting (Wang and Karimi, 2024). We view these applications as exciting avenues for future research.

Incorporating Human External Supervision. While DG-PRM constructs reward trees using LLMs, and although we have incorporated an automated validator to remove low-quality reward criteria, there remains the possibility of reward criteria that do not align with human expectations. This could lead to potential risks in optimized models. Therefore, it is crucial to introduce appropriate human supervision to refine and enhance the reward tree, and even to design and construct reward trees specific to certain domains.

Adaptation to Advanced Models. The primary strength of PRM lies in providing reliable and effective process signals for complex reasoning tasks. However, the most advanced models with deep reasoning capabilities, such as OpenAI’s o1 (OpenAI, 2024c) and o3 (OpenAI, 2025b), cannot be adapted due to their closed-source nature. A promising alternative is the DeepSeek R1 model (Guo et al., 2025). However, due to resource constraints, we have only utilized the distilled version of the model, rather than the full 671B R1 model. In the future, we plan to extend DG-PRM’s adaptability to more advanced models, further enhancing the effectiveness of generalized process supervision.

Ethics Statement

Compliance with Dataset Licenses. We strictly adhere to the licenses of the datasets used in our experiments. All datasets are in English, and we take care to ensure that our usage aligns with the intended use of each dataset. A detailed overview of the license information for each dataset is provided in Table 2.

Adherence to Model Usage Terms. Throughout the experimental process, we strictly follow the terms of use for the models. We comply with the usage guidelines set for models, including terms of service and API usage policies for commercial models. For open-source models, we adhere to the licenses and usage constraints outlined.

Data Privacy. Our method constructs reward trees using the output of LLMs without collecting personal information or sensitive data. We have thoroughly reviewed the prompts and data used in the experiments to ensure they do not contain any personally identifiable information or offensive content.

Data Annotation. During the experiments, we invited five annotators with a master’s degree or higher to label the process rewards of model outputs and assess DG-PRM outputs. One of the annotators, a PhD from the Chemistry Department, specifically handled ChemistryQA tasks. Compensation was provided to participants based on local hourly wage standards. All annotators were from the broader Pacific Rim region, with a balanced gender ratio, and we ensured that cultural preferences across different regions were accounted for in the evaluations, ensuring the validity of the results. The instructions provided to annotators are shown in Table 22 and Table 23.

Use of AI Assistants. In the evaluation process, we employed AI tools to assist with analyzing model outputs. Specifically, we utilized GitHub Copilot to assist in coding. We ensured that the use of AI tools followed submission guidelines and ethical standards.

Environmental Protection. Training and scaling during testing require significant computational power and resources. Efficient and accurate reward signals enable more efficient model training and allow for the elimination of lower-quality paths during inference, promoting the sustainable development of AI. This approach helps reduce carbon emissions and supports environmental protection.

Acknowledgment

This work was supported by the National Natural Science Foundation of China (No. U24B20181). The computations in this research were performed using the CFFF platform of Fudan University.

References

- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022a. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. 2022b. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.
- Yushi Bai, Jiahao Ying, Yixin Cao, Xin Lv, Yuze He, Xiaozhi Wang, Jifan Yu, Kaisheng Zeng, Yijia Xiao, Haozhe Lyu, Jiayin Zhang, Juanzi Li, and Lei Hou. 2023. Benchmarking foundation models with language-model-as-an-examiner. In *Advances in Neural Information Processing Systems*, volume 36, pages 78142–78167. Curran Associates, Inc.
- Meng Cao, Lei Shu, Lei Yu, Yun Zhu, Nevan Wichers, Yinxiao Liu, and Lei Meng. 2024. Enhancing reinforcement learning with dense rewards from language model critic. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 9119–9138, Miami, Florida, USA. Association for Computational Linguistics.
- Lu Chen, Rui Zheng, Binghai Wang, Senjie Jin, Caishuang Huang, Junjie Ye, Zhihao Zhang, Yuhao Zhou, Zhiheng Xi, Tao Gui, Qi Zhang, and Xuanjing Huang. 2024a. Improving discriminative capability of reward models in RLHF using contrastive learning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 15270–15283, Miami, Florida, USA. Association for Computational Linguistics.
- Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu, Mengfei Zhou, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. 2025. Do not think that much for $2+3=?$ on the overthinking of o1-like llms. *Preprint*, arXiv:2412.21187.
- Zhipeng Chen, Kun Zhou, Xin Zhao, Junchen Wan, Fuzheng Zhang, Di Zhang, and Ji-Rong Wen. 2024b. Improving large language models via fine-grained reinforcement learning with minimum editing constraint. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 5694–5711, Bangkok, Thailand. Association for Computational Linguistics.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. 2024. Chatbot arena: An open platform for evaluating llms by human preference. *Preprint*, arXiv:2403.04132.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *ArXiv*, abs/1803.05457.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *Preprint*, arXiv:2110.14168.
- Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang, Wendi Li, Bingxiang He, Yuchen Fan, Tianyu Yu, Qixin Xu, Weize Chen, Jiarui Yuan, Huayu Chen, Kaiyan Zhang, Xingtai Lv, Shuo Wang, Yuan Yao, Xu Han, Hao Peng, Yu Cheng, Zhiyuan Liu, Maosong Sun, Bowen Zhou, and Ning Ding. 2025. Process reinforcement through implicit rewards. *Preprint*, arXiv:2502.01456.
- Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. 2024. Safe RLHF: Safe reinforcement learning from human feedback. In *The Twelfth International Conference on Learning Representations*.
- Shihan Dou, Yan Liu, Haoxiang Jia, Enyu Zhou, Limao Xiong, Junjie Shan, Caishuang Huang, Xiao Wang, Xiaoran Fan, Zhiheng Xi, Yuhao Zhou, Tao Ji, Rui Zheng, Qi Zhang, Tao Gui, and Xuanjing Huang. 2024. StepCoder: Improving code generation with reinforcement learning from compiler feedback. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4571–4585, Bangkok, Thailand. Association for Computational Linguistics.
- Bofei Gao, Zefan Cai, Runxin Xu, Peiyi Wang, Ce Zheng, Runji Lin, Keming Lu, Dayiheng Liu, Chang Zhou, Wen Xiao, Junjie Hu, Tianyu Liu, and Baobao Chang. 2024. Llm critics help catch bugs in mathematics: Towards a better mathematical verifier with natural language feedback. *Preprint*, arXiv:2406.14024.
- Leo Gao, John Schulman, and Jacob Hilton. 2023a. Scaling laws for reward model overoptimization. In *International Conference on Machine Learning*, pages 10835–10866. PMLR.
- Mingqi Gao, Jie Ruan, Renliang Sun, Xunjian Yin, Shiping Yang, and Xiaojun Wan. 2023b. Human-like summarization evaluation with chatgpt. *Preprint*, arXiv:2304.02554.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361.
- Amelia Glaese, Nat McAleese, Maja Trębacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth Rauh,

- Laura Weidinger, Martin Chadwick, Phoebe Thacker, Lucy Campbell-Gillingham, Jonathan Uesato, Po-Sen Huang, Ramona Comanescu, Fan Yang, Abigail See, Sumanth Dathathri, Rory Greig, Charlie Chen, Doug Fritz, Jaume Sanchez Elias, Richard Green, Soňa Mokrá, Nicholas Fernando, Boxi Wu, Rachel Foley, Susannah Young, Iason Gabriel, William Isaac, John Mellor, Demis Hassabis, Koray Kavukcuoglu, Lisa Anne Hendricks, and Geoffrey Irving. 2022. [Improving alignment of dialogue agents via targeted human judgements](#). *Preprint*, arXiv:2209.14375.
- Xinyu Guan, Li Lina Zhang, Yifei Liu, Ning Shang, Youran Sun, Yi Zhu, Fan Yang, and Mao Yang. 2025. [rstar-math: Small llms can master math reasoning with self-evolved deep thinking](#). *Preprint*, arXiv:2501.04519.
- 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*.
- 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](#). *NeurIPS*.
- Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. 2024. [A survey of reinforcement learning from human feedback](#). *Preprint*, arXiv:2312.14925.
- Tushar Khot, Peter Clark, Michal Guerquin, Peter Jansen, and Ashish Sabharwal. 2020. [Qasc: A dataset for question answering via sentence composition](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, 05, pages 8082–8090.
- Minae Kwon, Sang Michael Xie, Kalesha Bullard, and Dorsa Sadigh. 2023. [Reward design with language models](#). In *The Eleventh International Conference on Learning Representations*.
- Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. 2024a. [Nv-embed: Improved techniques for training llms as generalist embedding models](#). *arXiv preprint arXiv:2405.17428*.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, and Sushant Prakash. 2024b. [Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback](#). *Preprint*, arXiv:2309.00267.
- Chaofan Li, MingHao Qin, Shitao Xiao, Jianlyu Chen, Kun Luo, Yingxia Shao, Defu Lian, and Zheng Liu. 2024a. [Making text embedders few-shot learners](#). *Preprint*, arXiv:2409.15700.
- Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhat-tacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, Kai Shu, Lu Cheng, and Huan Liu. 2024b. [From generation to judgment: Opportunities and challenges of llm-as-a-judge](#). *arXiv preprint arXiv: 2411.16594*.
- Junlong Li, Shichao Sun, Weizhe Yuan, Run-Ze Fan, hai zhao, and Pengfei Liu. 2024c. [Generative judge for evaluating alignment](#). In *The Twelfth International Conference on Learning Representations*.
- Wendi Li and Yixuan Li. 2024. [Process reward model with q-value rankings](#). *Preprint*, arXiv:2410.11287.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. [Let’s verify step by step](#). In *The Twelfth International Conference on Learning Representations*.
- Zhan Ling, Yunhao Fang, Xuanlin Li, Zhiao Huang, Mingu Lee, Roland Memisevic, and Hao Su. 2023. [Deductive verification of chain-of-thought reasoning](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Chris Yuhao Liu, Liang Zeng, Jiakai Liu, Rui Yan, Jujie He, Chaojie Wang, Shuicheng Yan, Yang Liu, and Yahui Zhou. 2024. [Skywork-reward: Bag of tricks for reward modeling in llms](#). *arXiv preprint arXiv:2410.18451*.
- Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Meiqi Guo, Harsh Lara, Yunxuan Li, Lei Shu, Yun Zhu, Lei Meng, Jiao Sun, and Abhinav Rastogi. 2024. [Improve mathematical reasoning in language models by automated process supervision](#). *Preprint*, arXiv:2406.06592.
- Dakota Mahan, Duy Van Phung, Rafael Rafailov, Chase Blagden, Nathan Lile, Louis Castricato, Jan-Philipp Fränken, Chelsea Finn, and Alon Albalak. 2024. [Generative reward models](#). *Preprint*, arXiv:2410.12832.
- Nat McAleese, Rai Michael Pokorny, Juan Felipe Ceron Uribe, Evgenia Nitishinskaya, Maja Trebacz, and Jan Leike. 2024. [Llm critics help catch llm bugs](#). *Preprint*, arXiv:2407.00215.
- Yu Meng, Mengzhou Xia, and Danqi Chen. 2024. [SimPO: Simple preference optimization with a reference-free reward](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Tong Mu, Alec Helyar, Johannes Heidecke, Joshua Achiam, Andrea Vallone, Ian D Kivlichan, Molly Lin, Alex Beutel, John Schulman, and Lilian Weng. 2024. [Rule based rewards for language model safety](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- OpenAI. 2024a. [Hello gpt-4o](#).
- OpenAI. 2024b. [New embedding models and api updates](#). Accessed: 2025-02-13.

- OpenAI. 2024c. [Openai o1-mini: Advancing cost-efficient reasoning.](#)
- OpenAI. 2025a. [Introducing deep research.](#)
- OpenAI. 2025b. [Openai o3-mini. Release.](#)
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback.](#) In *Advances in Neural Information Processing Systems*.
- Rafael Rafailov, Joey Hejna, Ryan Park, and Chelsea Finn. 2024. [From \$\\$r\$ to \$\\$q^*\$: Your language model is secretly a q-function.](#) In *First Conference on Language Modeling*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. [Direct preference optimization: Your language model is secretly a reward model.](#) In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Amrith Setlur, Chirag Nagpal, Adam Fisch, Xinyang Geng, Jacob Eisenstein, Rishabh Agarwal, Alekh Agarwal, Jonathan Berant, and Aviral Kumar. 2025. [Rewarding progress: Scaling automated process verifiers for LLM reasoning.](#) In *The Thirteenth International Conference on Learning Representations*.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. [Scaling llm test-time compute optimally can be more effective than scaling model parameters.](#) *Preprint*, arXiv:2408.03314.
- Mingyang Song, Zhaochen Su, Xiaoye Qu, Jiawei Zhou, and Yu Cheng. 2025. [Prmbench: A fine-grained and challenging benchmark for process-level reward models.](#) *arXiv preprint arXiv:2501.03124*.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. [Learning to summarize with human feedback.](#) In *Advances in Neural Information Processing Systems*, volume 33, pages 3008–3021. Curran Associates, Inc.
- Hao Sun, Thomas Pouplin, Nicolás Astorga, Tennison Liu, and Mihaela van der Schaar. 2024a. [Improving LLM generation with inverse and forward alignment: Reward modeling, prompting, fine-tuning, and inference-time optimization.](#) In *The First Workshop on System-2 Reasoning at Scale, NeurIPS'24*.
- Shengjie Sun, Runze Liu, Jiafei Lyu, Jing-Wen Yang, Liangpeng Zhang, and Xiu Li. 2024b. [A large language model-driven reward design framework via dynamic feedback for reinforcement learning.](#) *Preprint*, arXiv:2410.14660.
- Richard S Sutton. 2018. Reinforcement learning: An introduction. *A Bradford Book*.
- Qwen Team. 2024. [Qwq: Reflect deeply on the boundaries of the unknown.](#)
- Li Tianle, Chiang Wei-Lin, Frick Evan, Dunlap Lisa, Zhu Banghua, Joseph E. Gonzalez, and Stoica Ion. 2024. [From live data to high-quality benchmarks: The arena-hard pipeline.](#)
- 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.](#) *Preprint*, arXiv:2211.14275.
- Jonathan Uesato, Nate Kushman, Ramana Kumar, H. Francis Song, Noah Yamamoto Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. 2023. [Solving math word problems with process-based and outcome-based feedback.](#)
- Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. 2024a. [Math-shepherd: Verify and reinforce LLMs step-by-step without human annotations.](#) In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9426–9439, Bangkok, Thailand. Association for Computational Linguistics.
- Yang Wang and Hassan A. Karimi. 2024. [Exploring large language models for climate forecasting.](#) *Preprint*, arXiv:2411.13724.
- Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy J. Zhang, Makesh Narsimhan Sreedhar, and Oleksii Kuchaiev. 2024b. [Helpsteer2: Open-source dataset for training top-performing reward models.](#) *Preprint*, arXiv:2406.08673.
- Zihan Wang, Yunxuan Li, Yuexin Wu, Liangchen Luo, Le Hou, Hongkun Yu, and Jingbo Shang. 2024c. [Multi-step problem solving through a verifier: An empirical analysis on model-induced process supervision.](#) In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7309–7319, Miami, Florida, USA. Association for Computational Linguistics.
- Joe H Ward Jr. 1963. Hierarchical grouping to optimize an objective function. *Journal of the American statistical association*, 58(301):236–244.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. [Chain of thought prompting elicits reasoning in large language models.](#) In *Advances in Neural Information Processing Systems*.
- Zhuoyu Wei, Wei Ji, Xiubo Geng, Yining Chen, Baihua Chen, Tao Qin, and Daxin Jiang. 2021. [Chemistry{qa}: A complex question answering dataset from chemistry.](#)

- Jiaxin Wen, Ruiqi Zhong, Akbir Khan, Ethan Perez, Jacob Steinhardt, Minlie Huang, Samuel R. Bowman, He He, and Shi Feng. 2025. [Language models learn to mislead humans via RLHF](#). In *The Thirteenth International Conference on Learning Representations*.
- Lilian Weng. 2024. [Reward hacking in reinforcement learning](#). *lilianweng.github.io*.
- Jianchang Wu, Luca Torresi, ManMan Hu, Patrick Reiser, Jiyun Zhang, Juan S. Rocha-Ortiz, Luyao Wang, Zhiqiang Xie, Kaicheng Zhang, Byung wook Park, Anastasia Barabash, Yicheng Zhao, Junsheng Luo, Yunuo Wang, Larry Lüer, Lin-Long Deng, Jens A. Hauch, Dirk M. Guldi, M. Eugenia Pérez-Ojeda, Sang Il Seok, Pascal Friederich, and Christoph J. Brabec. 2024. [Inverse design workflow discovers hole-transport materials tailored for perovskite solar cells](#). *Science*, 386(6727):1256–1264.
- Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024. [Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation](#). In *ICML*.
- 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, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. 2024b. [Qwen2.5-math technical report: Toward mathematical expert model via self-improvement](#). *arXiv preprint arXiv:2409.12122*.
- Lifan Yuan, Wendi Li, Huayu Chen, Ganqu Cui, Ning Ding, Kaiyan Zhang, Bowen Zhou, Zhiyuan Liu, and Hao Peng. 2024. [Free process rewards without process labels](#). *Preprint*, arXiv:2412.01981.
- Tian Zhang, Raghu Ramakrishnan, and Miron Livny. 1997. [Birch: A new data clustering algorithm and its applications](#). *Data mining and knowledge discovery*, 1:141–182.
- 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](#). *Preprint*, arXiv:2501.07301.
- Chujie Zheng, Zhenru Zhang, Beichen Zhang, Runji Lin, Keming Lu, Bowen Yu, Dayiheng Liu, Jingren Zhou, and Junyang Lin. 2024a. [Processbench: Identifying process errors in mathematical reasoning](#). *Preprint*, arXiv:2412.06559.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). *Preprint*, arXiv:2306.05685.
- Yizhen Zheng, Huan Yee Koh, Maddie Yang, Li Li, Lauren T. May, Geoffrey I. Webb, Shirui Pan, and George Church. 2024b. [Large language models in drug discovery and development: From disease mechanisms to clinical trials](#). *Preprint*, arXiv:2409.04481.
- Shuang Zhou, Zidu Xu, Mian Zhang, Chunpu Xu, Yawen Guo, Zaifu Zhan, Sirui Ding, Jiashuo Wang, Kaishuai Xu, Yi Fang, Liqiao Xia, Jeremy Yeung, Daochen Zha, Genevieve B. Melton, Mingquan Lin, and Rui Zhang. 2024. [Large language models for disease diagnosis: A scoping review](#). *Preprint*, arXiv:2409.00097.
- Lianghui Zhu, Xinggang Wang, and Xinlong Wang. 2024. [JudgeLM : Fine-tuned large language models are scalable judges](#).

A Dataset Details

In our experiments, we selected seven datasets encompassing a wide variety of task types that require intricate and complex reasoning by the models. Detailed information on the sample sizes, sources, and licenses of these datasets is provided in Table 1.

- **PRMBENCH** is evaluated using the PRM-Score, which is the average of the F1 Score and Negative F1 Score, with an emphasis on the steps where errors occur. Table 3 provides a detailed list of the abbreviations for each category, their corresponding full names, the evaluation objectives, and the number of instances within each class. An example is shown in Table 8.
- **MT-Bench** (Zheng et al., 2023) is evaluated on a 1-10 scale. To ensure comprehensive evaluation and realistic scenarios, we employed a multi-turn setup with evaluation scores output by GPT-4o (OpenAI, 2024a). We also report the win rates against the GPT-4o model. Table 4 presents the various types and their corresponding quantities in MT-Bench, and Table 13 illustrates an example from MT-Bench.
- **Arena-Hard** (Tianle et al., 2024) uses win rates against the GPT-4-0314 output as an evaluation metric. We used the official output results from [Arena-Hard’s official repository](#), utilizing GPT-4o (OpenAI, 2024a) as the judge model. Arena-Hard consists of 250 diverse scenarios. An example is shown in Table 14.

DATASET	TASK	ANSWER FORMAT	# TRAIN.	# TEST.	LICENSE
PRMBENCH (Song et al., 2025)	PRM	Incorrect Position	-	6,216	Apache license 2.0
MT-Bench (Zheng et al., 2023)	General	Rating	-	80	Apache license 2.0
Arena-Hard (Tianle et al., 2024)	General	Rating	-	500	Apache license 2.0
QASC (Khot et al., 2020)	Science	Multi-choice	8,134	926	Apache license 2.0
ChemistryQA (Wei et al., 2021)	Science	Text	2,721	392	CC0-1.0
StrategyQA (Geva et al., 2021)	Commonsense	T/F	2,061	229	MIT license
ARC-c (Clark et al., 2018)	Commonsense	Multi-choice	1,119	299	CC BY-SA 4.0

Table 2: Detailed description of the datasets used in the experiments. # TRAIN. represents the number of training samples, and # TEST. indicates the number of samples used for evaluation.

Full Name	Evaluation Subject	Number
Non-Redundancy (NR.)	Simplicity	758
Non-Circular Logic (NCL.)	Simplicity	758
Empirical Soundness (ES.)	Soundness	757
Step Consistency (SC.)	Soundness	758
Domain Consistency (DC.)	Soundness	757
Confidence Invariance (CI.)	Soundness	757
Prerequisite Sensitivity (PS.)	Sensitivity	756
Deception Resistance (DR.)	Sensitivity	750
Multi-Solution Consistency (MS.)	Sensitivity	165

Table 3: Statistics of classes in MT-Bench.

- **QASC** (Khot et al., 2020) is a multiple-choice science dataset that provides multiple fact explanations for each answer, aiding LLMs in evaluating the causes of errors from the perspective of facts. A simple example is shown in Table 9.
- **ChemistryQA** (Wei et al., 2021) is a chemistry dataset collected by Socratic, covering over 200 topics. It provides the necessary knowledge, conditions, and detailed solution steps. A sample is shown in Table 10.
- **StrategyQA** (Geva et al., 2021) challenges models to solve implicit reasoning problems strategically. It provides question breakdowns and corresponding facts, which help identify the points of failure. A simple example is shown in Table 11.
- **ARC-c** (Clark et al., 2018) is a challenge subset of the ARC dataset, assessing a model’s fundamental reasoning abilities across various fact types, such as Basic Facts & Properties, and Processes & Causal.

Since PRMBENCH, MT-Bench, and Arena-Hard do not provide training datasets suitable for constructing reward trees, we leverage the MATH dataset (Hendrycks et al., 2021) to train PRMBENCH’s reward tree. For MT-Bench and Arena-Hard, we select 5,000 samples from the LMSYS-Human-Preference-55k (Chiang et al., 2024) dataset to build the reward trees.

Task Category	Evaluation Focus	Number of Samples
Writing	Text Generation	10
Roleplay	Interaction	10
Reasoning	Logical Analysis	10
Math	Mathematical Problem Solving	10
Coding	Programming	10
Extraction	Information Retrieval	10
STEM	Scientific Knowledge	10
Humanities	Cultural Understanding	10

Table 4: Statistics of categories in MT-Bench.

B Experiment Details

B.1 Implementation Details

In the process, we segment steps using newline characters or explicit labels such as “Step1” and “Step2.” For each step, we construct positive and negative label pairs based on the reference answers y or evaluation criteria \mathcal{C} provided in the dataset, such as MT-Bench (Zheng et al., 2023). We use the corresponding judge prompt and the reference answer y from the dataset to build these pairs.

We utilize GPT-4o to construct the judge model \mathcal{J} , generating appropriate reward criteria for erroneous steps. The BAAI/bge-en-icl (Li et al., 2024a) embedding model is employed to obtain embedding vectors v , where the vector space has a dimension $d = 4096$. Additionally, in our experimental analysis, we incorporate the text-embedding-3-large (OpenAI, 2024b) and nvidia/NV-Embed-v2 (Lee et al., 2024a) models.

The BIRCH algorithm (Zhang et al., 1997) is used for hierarchical clustering, with the Birch implementation from sklearn.cluster. In the process of merging reward criteria, we retain criteria with longer lengths. Following the clustering, GPT-4o is used to summarize the upper-level coarse-grained reward criteria r^{parent} . This step proves crucial for subsequent reward allocation, as using higher-quality models for summarization yields more accurate reward signals.

The same model used to generate the coarse-grained reward criteria is employed for reward al-

location, determining whether fine-grained reward criteria r^{child} should be provided. To ensure the accuracy of evaluation, we sample multiple times to confirm the appropriateness of selected reward criteria. In our experiments, we set the sampling frequency to 5, retaining coarse-grained reward criteria that appear more than 3 times. Typically, the model’s judgments are consistent.

During scoring, previously generated steps, evaluation information, and scores are concatenated at the front for reference. By default, we use GPT-4o as the scorer \mathcal{S} . We also analyze the evaluation results using the model’s own evaluation and the o3-mini evaluation results, as detailed in D.1.

If corresponding positive or negative feedback cannot be found, additional sampling is performed until the maximum number of attempts is reached. Table 18 provides an example with selected fine-grained reward criteria placed in the prompt. In practice, only coarse-grained reward criteria can be included due to context limitations, so we carefully select fine-grained reward criteria based on matching. In Table 19, Table 20, and Table 21, we present several examples of coarse-grained and fine-grained process reward criteria.

For training, we use 8 interconnected H200 GPUs with a batch size of 32, a learning rate of $5e-6$, and a DPO-beta of 0.1. Other settings are based on the default parameters of HuggingFace’s DPOTrainer. Due to the temporary unavailability of the DeepSeek-R1 API, we deploy the full 671B version using ollama on an 8 interconnected A100 setup with INT4 precision. The same approach is used for the other Distill models.

During generation, we adjust the temperature parameter γ according to the task and model. Specifically, we find that for R1-Distill-Qwen-1.5B and R1-Distill-Qwen-7B, higher temperatures lead to very long sequences and many repeated solutions, as observed by Chen et al. (2025) in the QwQ (Team, 2024) model. Thus, we set $\gamma \in [0.5, 0.6]$ and filter out sequences exceeding 4096 tokens. For R1-Distill-Qwen-14B and R1-Distill-Qwen-32B, we find that $\gamma \in [0.7, 0.8]$ yields more satisfactory results on MT-Bench and Areal-Hard, while a lower temperature works better for QASC, StrategyQA, and ARC-c tasks, so we set $\gamma \in [0.6, 0.7]$. For the MT-Bench and Areal-Hard datasets, we compute the average score of five output results. For other datasets, we calculate the final results based on the corresponding metrics.

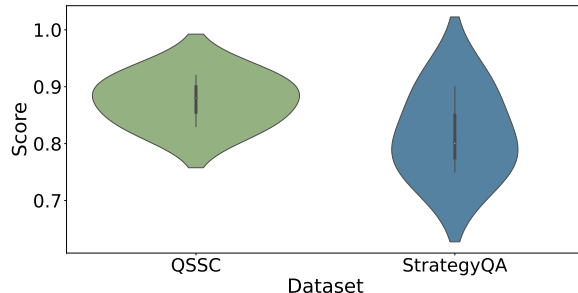


Figure 7: Consistency between automated validator and human evaluation

B.2 Automated Validator

During the construction of the reward tree \mathcal{T} , an automated validator plays a crucial role in ensuring the high quality of the reward criteria. This validator can be applied to the reward tree construction for any model, filtering out low-quality reward criteria. We use GPT-4o (OpenAI, 2024a) as the automated validator, which evaluates the reward criteria through prompts and provides one of three assessment results: Good, Ordinary and Bad. Detailed criteria for these evaluations are outlined, and Table 15 provides a general example.

To assess the consistency between the automated validator and human evaluations, we randomly selected 100 outputs from the automated validator and compared them to the results of human evaluations. We calculated the proportion of matches between the automated outputs and the labels assigned by three human evaluators, with the results shown in Figure 7. Our analysis reveals that, on both the QASC and StrategyQA datasets, the automated validator exhibits high consistency with human evaluations.

B.3 Baseline Implementation

For the ORM, we follow the approach of Liu et al. (2024) and train the ORM based on the Qwen-2.5-7B model (Yang et al., 2024a). We employ AdamW as the optimizer, with a batch size of 16 and a learning rate of $2e-6$. In the process of building Critic Models, we prompt the LLM to perform step-by-step evaluation. Specifically, the Critic Direct directly outputs whether each step is correct, while the Critic CoT provides a detailed analysis of the output steps before giving the final result. Each step is evaluated with `[[correct]]` or `[[wrong]]`, and the final judge model extracts any erroneous steps. If all steps are correct, the output is `[]`. Tables 16 and 17 show an example of the Critic Direct and Critic CoT.

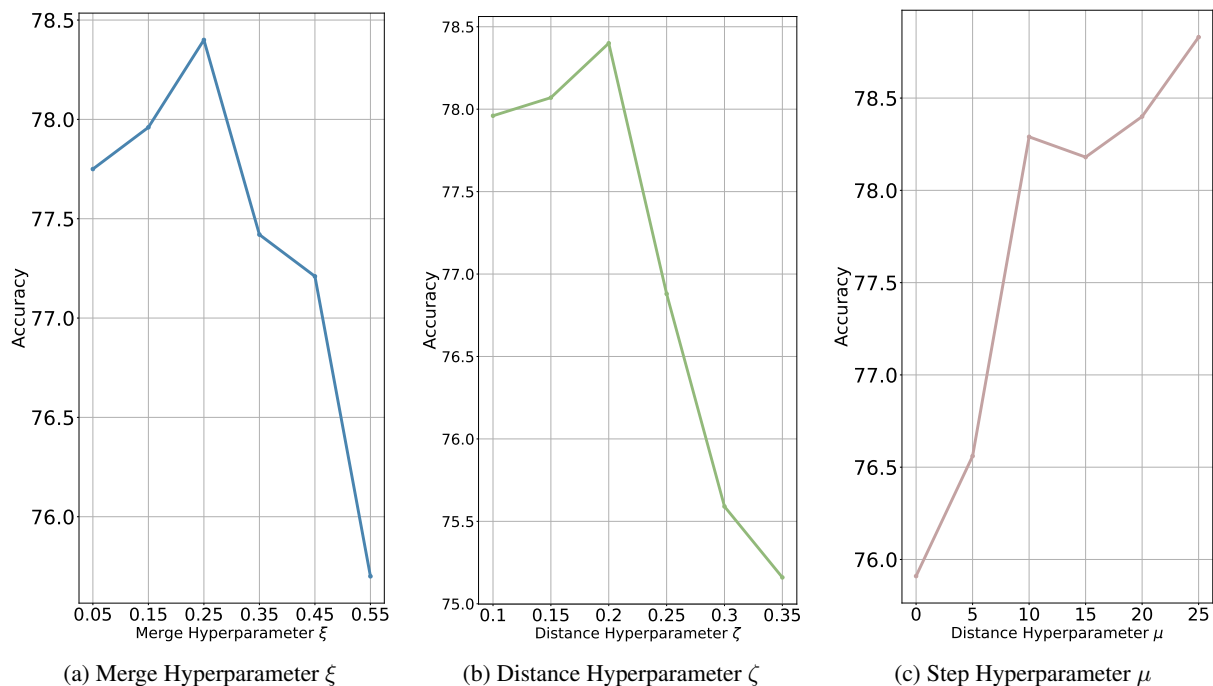


Figure 8: Ablation analysis of the hyperparameters on the QASC dataset, using R1-Distill-Qwen-7B as the backbone.

C Hyperparameter Analysis

We conducted an ablation study on the hyperparameters used in our experiments on the QASC dataset, as shown in Figure 8.

Merge Hyperparameter ξ The hyperparameter ξ controls the merging of similar reward criteria. As shown in Figure 8a, a larger value of ξ leads to the merging of related criteria, and in an extreme case, all r^{child} under r^{parent} are merged, resulting in a coarse-grained process reward. This has a significant impact on performance. On the other hand, a smaller value of ξ leads to a large number of redundant criteria in the reward tree \mathcal{T} , which increases noise and causes performance degradation. It also results in redundant computations in \mathcal{S} , leading to unnecessary overhead. Therefore, we set $\xi = 0.25$ as a suitable threshold.

Distance Hyperparameter ζ The hyperparameter ζ controls the selection of fine-grained process reward criteria, as illustrated in Figure 8b. A larger value of ζ results in the inclusion of numerous irrelevant reward criteria, leading to a decrease in overall performance. This occurs because an excessive number of criteria makes it difficult for DG-PRM to select the corresponding positive and negative samples. Consequently, it is essential to limit the criteria within an optimal range. On the other hand, a smaller value of ζ may fail to match the appropriate fine-grained process reward criteria,

similarly impairing performance.

Step Hyperparameter μ The hyperparameter μ controls the number of steps within the reward criteria selection \mathcal{R} and scoring \mathcal{S} that can be referenced. We observe that as the number of steps increases, performance continues to improve. Therefore, providing a larger number of steps is beneficial for the accuracy of DG-PRM selection and scoring. However, considering the constraints of the model’s context window, we do not set μ above 20. In the ChemistryQA scenario, the reasoning process is more complex, which can exceed the model’s context window. Additionally, we observe that performance gains gradually diminish. Therefore, considering the cost overhead, we set $\mu = 20$.

D Analysis and Discussion

D.1 Judge Model

In Figure 9, we analyze the impact of different models as \mathcal{R} and \mathcal{S} on performance using the MT-Bench dataset. We utilize radar charts to display the variation in ratings across each category. Our observations reveal that employing a better judge model can significantly enhance performance in writing, reasoning, and coding tasks, indicating that DG-PRM can continuously benefit from a superior judge model. Furthermore, we find that DeepSeek-R1 exhibits performance comparable to o3-mini, even outperforming it in coding and STEM tasks. This demonstrates that DG-PRM is also applicable

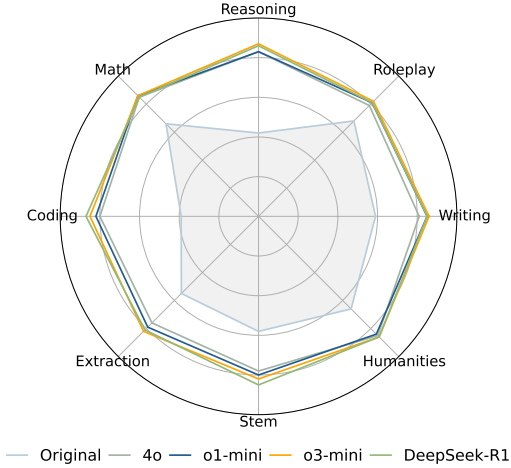


Figure 9: Impact of the judge model on ratings in the MT-Bench dataset, using R1-Distill-Qwen-7B as the backbone.

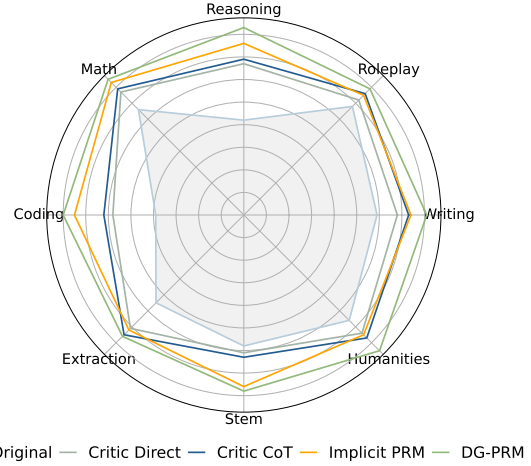


Figure 10: Performance comparison of different PRM methods on the MT-Bench dataset, using R1-Distill-Qwen-7B as the backbone.

Embedding Model	QASC	ChemistryQA
BAAI/bge-en-icl	78.40	87.50
text-embedding-3-large	79.04	88.01
nvidia/NV-Embed-v2	78.29	87.24

Table 5: Analysis of embedding models.

to advanced open-source models.

D.2 Implicit PRM

In Figure 10, we compare the performance of different methods on the MT-Bench dataset. We include Implicit PRM (Rafailov et al., 2024), which is trained using the LMSYS-Human-Preference-55k (Chiang et al., 2024) dataset based on the official implementation. We regenerate the responses using R1-Distill-Qwen-7B and construct the chosen and rejected pairs based on the scoring from GPT-4o. We observe that Implicit PRM effectively improves performance, particularly in reasoning and coding, suggesting that Implicit PRM can effectively model tasks with well-defined objectives. However, DG-PRM demonstrates more substantial improvements, offering a more comprehensive enhancement of model performance, such as its exceptional results in writing and humanities tasks. Furthermore, DG-PRM is more interpretable, providing a clear explanation of the advantages of positive samples over negative samples, making the optimization objectives easier to understand.

D.3 Embedding Model

Table 5 analyzes the impact of different embedding models on performance. The embedding model plays a crucial role in the construction of the re-

Embedding Model	QASC	ChemistryQA
BRICH	78.40	87.50
Agglomerative	78.83	86.73
Divisive	79.48	88.26

Table 6: Analysis of hierarchical clustering algorithm.

ward tree and the selection of fine-grained process reward criteria. We find that DG-PRM demonstrates robustness in terms of embedding model selection, achieving satisfactory performance even with open-source models. Therefore, we use the BAAI/bge-en-icl model as the embedding function \mathcal{V} in our experiments.

D.4 Hierarchical Clustering

In Table 6, we examine the effect of different hierarchical clustering methods on performance. We selected two approaches: agglomerative and divisive clustering, using Ward’s method (Ward Jr, 1963) to define cluster distances. Our findings reveal that divisive clustering yields better performance. However, given the high computational cost of divisive clustering, the incremental updates offered by the BRICH algorithm (Zhang et al., 1997) significantly reduce this overhead, making the addition and removal of reward criteria more efficient and convenient. Therefore, we selected the BRICH algorithm as the clustering method to obtain the reward tree.

D.5 General Reward Tree \mathcal{T}

To further evaluate the generalizability of DG-PRM, we construct a unified reward tree encompassing a rich set of criteria from the MATH, QASC, Chemistry, StrategyQA, and ARC-c train-

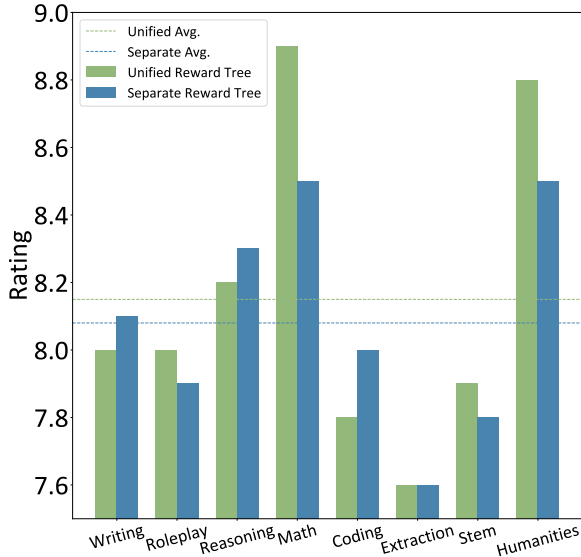


Figure 11: Comparison of ratings on MT-Bench using a unified reward tree versus separate reward trees.

ing datasets. Figure 11 presents the experimental results using this unified reward tree on MT-Bench. Compared to a reward tree constructed separately for MT-Bench, the unified reward tree demonstrates better performance on humanities and math tasks. We observe that the model incorporates a broader range of factual evaluations during scoring. These criteria likely stem from the diverse mathematical perspectives in MATH and fact-related assessments in StrategyQA. The dashed line represents the average performance of the unified and separate reward trees. Our findings show that the unified reward tree achieves a higher overall score, highlighting the exceptional scalability of DG-PRM.

D.6 Human Consensus

To further analyze the rationality of DG-PRM’s reward criteria allocation and output reward scores, we randomly select 100 samples, including the problem, answer, assigned criteria, and final reward score. These are independently evaluated by three assessors, who judge the appropriateness of the allocated criteria and output score, labeling each sample as Good, Ordinary, or Bad. The results are shown in Figure 12. We observe that the median proportion of Good labels is close to 80%, while the proportion of Bad labels is under 10%, indicating that the majority of assessors find DG-PRM’s reward allocation and output scores to be reasonable, demonstrating the rationality of DG-PRM’s reward criteria construction and score output.

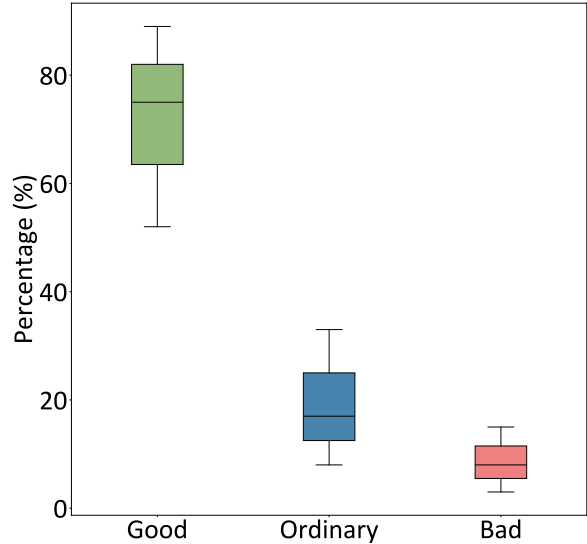


Figure 12: Human evaluation of DG-PRM output, including the selection of reward criteria and scores.

Optimization Algorithm	Single-Turn	Multi-Turn
Baseline	6.01	5.66
DPO (Rafailov et al., 2023)	8.67	8.09
CPO (Xu et al., 2024)	8.58	7.87
SimPO (Meng et al., 2024)	8.62	7.95

Table 7: Performance analysis of different optimization algorithms on MT-Bench.

D.7 Optimization Algorithm

In Table 7, we analyze the performance of DG-PRM using various optimization algorithms. Experimental results on MT-Bench show that the process rewards constructed by DG-PRM are effective across multiple optimization algorithms. Significant performance improvements are observed in both Single-Turn and Multi-Turn scenarios compared to the baseline. Among the algorithms tested, DPO achieves the most notable performance gain, thus we utilize DPO to optimize our policy model in the experiments.

Visions of Tomorrow

In recent years, we have witnessed a rapid expansion of model capabilities, with AI increasingly integrating into human scientific research (OpenAI, 2025a). This necessitates the development of a generally applicable Process Reward Model (PRM). In the short term, PRM will focus on aligning with the preferences of human scientists and our understanding of the objective world, preventing erroneous search paths and guiding the model toward more meaningful exploration. As model capabilities continue to evolve, PRM will steer the model

towards human-desired goals, such as environmental sustainability and disease treatment, rather than facilitating the creation of harmful products in explorations. We believe that a generalizable PRM will play a pivotal role in advancing AI in a safer and more reliable direction.

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: $x^3 - 2x + 1 = 2$.

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

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.

Table 8: PRMBENCH samples.

Question: Climate is generally described in terms of what?

Answer Choice: (A) sand (B) occurs over a wide range (C) forests (D) Global warming (E) rapid changes occur (F) local weather conditions (G) measure of motion (H) city life

Fact1: Climate is generally described in terms of temperature and moisture.

Fact2: Fire behavior is driven by local weather conditions such as winds, temperature and moisture.

answerKey: F

Table 9: QASC samples.

Question: How many moles of phosphorus trichloride would contain 3.35 times 10^{24} molecules of phosphorus trichloride?

Conditions: Number [OF] phosphorus trichloride molecules [=]pu3.3510²⁴

TargetVariable: Mole [OF] phosphorus trichloride [IN] moles

Answer: 5.56 moles

Table 10: ChemistryQA samples.

Question: Did Aristotle use a laptop?

Answer: No

Table 11: StrategyQA samples.

Question: Juan and LaKeisha roll a few objects down a ramp. They want to see which object rolls the farthest. What should they do so they can repeat their investigation?

Answer Choice: (A) Put the objects in groups. (B) Change the height of the ramp. (C) Choose different objects to roll. (D) Record the details of the investigation.

Answer: D

Table 12: ARC-c samples.

Question 1: Compose an engaging travel blog post about a recent trip to Hawaii, highlighting cultural experiences and must-see attractions.

Answer **A Glimpse into the Enchanted Land of Hawaii: A Cultural and Adventurous Journey**

Our recent trip to Hawaii was an absolute blast, a journey that left us thoroughly charmed by its beauty, culture, and unique spirit. As we explored this archipelago of islands, each destination offered a taste of its rich heritage and stunning landscapes.

One experience that stood out for me was the traditional Luau at Waikiki Beach Park. What made it truly unforgettable was the opportunity to step into a cultural immersion, dancing to the rhythm of hula while sharing stories with locals. It was like being transported back in time to a vibrant Hawaiian village, where music and hospitality were the foundation of life.

Our journey began in Honolulu, where we strolled along Waikiki Beach, soaking in the warm sun and vibrant ocean views. This iconic location isn't just a beach; it's a symbol of Hawaii's blend of tradition and modernity, perfect for capturing photos with the iconic "Hawaii Five-O" lighthouse in the background.

Next up was Pearl Harbor, where history took center stage. The experience was deeply moving as we witnessed the power of remembrance at this historic site. It reminded us how crucial it is to honor our heroes and preserve their legacy.

We then ventured to Hanauma Bay, a destination steeped in nature's beauty. Scattered with fronds and lined by sea grasses, this bay offered a serene escape, allowing us to truly disconnect from the hustle and bustle of city life.

Finally, we visited the Ironman Training Center at Waikoloa State Park, challenging ourselves physically and mentally. The rigorous workout left us exhilarated, pushing our limits while fostering camaraderie with fellow participants.

Reflecting on this trip, I can say that Hawaii's culture is as vibrant as its landscapes. Every experience was a lesson in appreciation for life's beauty and the richness of diverse traditions. Our time there will forever be cherished, not just for the sights and sounds, but for the profound impact it had on our hearts.

Hawaii's magic continues to inspire us, and we can't wait to return. If you're ever in the mood for an adventurous adventure with a touch of culture, Hawaii is the perfect destination!

Question 2: Rewrite your previous response. Start every sentence with the letter A.

Answer:

A Clever Rewrite: A Glimpse into the Enchanted Land of Hawaii

A recent trip to Hawaii was an absolute adventure, a journey that left us thoroughly charmed by its beauty, culture, and unique spirit. A highlight of our trip was experiencing traditional Hawaiian culture at a luau. A visit to Waikiki Beach Park offered a chance to dance to the rhythm of hula while sharing stories with locals—a true cultural immersion.

A beginning in Honolulu saw us strolling along Waikiki Beach, soaking in the warm sun and vibrant ocean views—a quintessential location for capturing photos with the iconic "Hawaii Five-O" lighthouse. A stop at Pearl Harbor provided a powerful reminder of history's significance. A visit to Hanauma Bay offered a serene escape, lined by sea grasses and scattered with fronds.

A challenge we faced during our adventure was finding words that started with "A," but I managed to overcome it by rephrasing sentences creatively. For example, instead of saying "We arrived in Honolulu," I wrote, "Arriving in Honolulu was an exciting beginning." Each sentence now starts with an "A," making the blog post a bit unconventional yet engaging. Despite this unusual formatting, the overall message remains clear and appealing—Hawaii's culture is as vibrant as its landscapes. Every experience was a lesson in appreciating life's beauty and the richness of diverse traditions. Our time there will forever be cherished, not just for the sights and sounds, but for the profound impact it had on our hearts.

A trip to Hawaii is an absolute must for anyone seeking an adventurous adventure with a touch of culture. A destination like Hawaii is perfect! We can't wait to return.

If you're in the mood for an adventurous adventure with a touch of culture, A trip to Hawaii could be just what you need. A wonderful experience that will stay with you forever!

Judge

Please act as an impartial judge and evaluate the quality of the response provided by an AI assistant to the user question displayed below. Your evaluation should consider factors such as the helpfulness, relevance, accuracy, depth, creativity, and level of detail of the response. Your evaluation should focus on the assistant's answer to the second user question. Begin your evaluation by providing a short explanation. Be as objective as possible. After providing your explanation, you must rate the response on a scale of 1 to 10 by strictly following this format: "[rating]", for example: "Rating: [5]".

Table 13: MT-Bench samples.

Question: Help me draft a research introduction of this topic "Data-Driven Insights into the Impact of Climate and Soil Conditions on Durian Floral Induction"

Answer: ****Introduction to the Impact of Climate and Soil Conditions on Durian Floral Induction****

Durian, a fragrant and sweet tropical fruit, is a cropping staple in various regions of Southeast Asia, particularly in countries like Malaysia and Thailand. As a major agricultural commodity, its economic significance is substantial, contributing significantly to local economies. However, durian cultivation faces challenges, especially in the face of climate change, which can impact agricultural productivity globally.

Climate change poses particular threats to durian cultivation, affecting factors such as temperature, humidity, and rainfall, all of which are critical for optimal flowering. Floral induction, a horticultural practice aimed at promoting flowering in plants to enhance fruit production, is essential for durian cultivation. Similarly, soil conditions, including fertility, pH levels, and organic matter, play a pivotal role in influencing the flowering process. Some durian varieties are particularly sensitive to these environmental factors, necessitating a deeper understanding of their interplay.

This study seeks to explore how climate and soil conditions specifically impact the floral induction of durian, with the aim of providing insights that can enhance sustainable agricultural practices. By identifying optimal conditions, the research contributes to improving durian yields and efficiency, offering valuable information for farmers and policymakers. This investigation not only addresses existing gaps in current knowledge but also provides a foundation for developing more resilient durian varieties and farming practices, ensuring continued productivity in the face of environmental challenges.

Judge

Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants to the user prompt displayed below. You will be given assistant A's answer and assistant B's answer. Your job is to evaluate which assistant's answer is better.

Begin your evaluation by generating your own answer to the prompt. You must provide your answers before judging any answers.

When evaluating the assistants' answers, compare both assistants' answers with your answer. You must identify and correct any mistakes or inaccurate information.

Then consider if the assistant's answers are helpful, relevant, and concise. Helpful means the answer correctly responds to the prompt or follows the instructions. Note when user prompt has any ambiguity or more than one interpretation, it is more helpful and appropriate to ask for clarifications or more information from the user than providing an answer based on assumptions. Relevant means all parts of the response closely connect or are appropriate to what is being asked. Concise means the response is clear and not verbose or excessive.

Then consider the creativity and novelty of the assistant's answers when needed. Finally, identify any missing important information in the assistants' answers that would be beneficial to include when responding to the user prompt.

After providing your explanation, you must output only one of the following choices as your final verdict with a label:

1. Assistant A is significantly better: [[A»B]]
2. Assistant A is slightly better: [[A>B]]
3. Tie, relatively the same: [[A=B]]
4. Assistant B is slightly better: [[B>A]]
5. Assistant B is significantly better: [[B»A]]

Example output: "My final verdict is tie: [[A=B]]".

Table 14: Arena-Hard samples.

Prompt: You are given a reward criterion used to evaluate a task or output from a model. Your task is to categorize the quality of this reward criterion into one of the following three levels: Good, Ordinary, or Bad. Please ensure your categorization is based on the following specific standards for each level.

Good Reward Criterion:

- The criterion is highly relevant to the task, focusing on evaluating the most important and central aspects of the response. It aligns well with the goals and requirements of the task.
- The criterion encourages the generation of high-quality, focused, and precise outputs. It promotes clarity, accuracy, and thoroughness, ensuring that key points are well-addressed.
- The criterion strikes a balance between being neither too broad nor too narrow, ensuring it is applicable across a range of responses while not overlooking key details. It facilitates easy application and optimal model adjustment.
- The criterion is practical, straightforward, and well-defined, allowing for consistent and reproducible evaluations. It minimizes ambiguity in interpretation.

Output: [[Good]]

Ordinary Reward Criterion:

- The criterion is somewhat specific but lacks the necessary depth or precision to effectively evaluate all important aspects of the response. It may miss key considerations or leave room for misinterpretation.
- The criterion aligns with the task's goal but may fail to account for subtle nuances or complex requirements. It could be improved by incorporating more detail or covering a broader scope.
- While the criterion provides guidance, it might not lead to optimal model adjustments due to its vagueness or complexity. It may be harder to apply consistently.
- The criterion is somewhat practical but lacks clarity or specificity in certain areas, potentially leading to inconsistent evaluation or suboptimal feedback.

Output: [[Ordinary]]

Bad Reward Criterion:

- The criterion is irrelevant to the task or focuses on factors that do not contribute meaningfully to evaluating the response, such as unrelated performance metrics or non-essential features.
- The criterion is overly complex, ambiguous, or misleading, making it difficult to understand or apply. It may introduce unnecessary confusion or fail to directly align with the task.
- The criterion encourages the wrong kind of response, such as penalizing desirable features (e.g., creativity, humor, or engagement) or evaluating irrelevant aspects of the output that do not pertain to the task's goal.
- The criterion is impractical or overly rigid, leading to evaluations that are either too subjective or too arbitrary. It may discourage meaningful feedback or improvements.

Output: [[Bad]]

Table 15: Automated Validator Prompt.

Prompt: Please evaluate the reasoning process below and directly output the erroneous step. If the step is correct, please output [[]]. You may refer to the following example.

Question: What type of water formation is formed by clouds?

Answer Choice: (A) pearls (B) streams (C) shells (D) diamonds (E) rain (F) beads (G) cooled (H) liquid

Step 1: Understanding the question.

The question asks, "What type of water formation is formed by clouds?" We need to determine the correct type of water formation typically associated with clouds.

Step 2: Analyzing the options.

The options provided are: (A) pearls (B) streams (C) shells (D) diamonds (E) rain (F) beads (G) cooled (H) liquid

Step 3: Identify the key concept.

Clouds are made up of tiny water droplets or ice crystals that have condensed from vapor in the atmosphere. These droplets can eventually combine and fall to the ground as precipitation.

Step 4: Eliminate irrelevant options.

Pearls, streams, shells, diamonds, beads, and cooled do not seem to be directly related to water formation associated with clouds.

Step 5: Focus on the most plausible option.

"Rain" is the most likely type of water formation that clouds produce, as clouds release water in the form of rain when the droplets or ice crystals become heavy enough to fall.

Final answer: [[E]]

Step 1: [[correct]]

Step 2: [[correct]]

Step 3: [[correct]]

Step 4: [[wrong]]

Step 5: [[correct]]

Judge: [[4]]

Question: Where do beads of water come from?

Answer Choice: (A) Too much water (B) underground systems (C) When the water is too cold (D) Water spills (E) Vapor turning into a liquid (F) Warm air moving into cold air (G) At the peak of a mountain (H) To another location like underground

Step 1: Understanding the question.

The question asks where beads of water come from. Beads of water are typically small droplets that form when water vapor condenses.

Step 2: Analyzing the options.

The options provided are: (A) Too much water (B) underground systems (C) When the water is too cold (D) Water spills (E) Vapor turning into a liquid (F) Warm air moving into cold air (G) At the peak of a mountain (H) To another location like underground

Step 3: Focus on key concepts.

Beads of water typically form when water vapor cools and condenses into liquid droplets.

Step 4: Eliminate irrelevant options.

(A) "Too much water" is not a correct explanation for beads of water.

(B) "Underground systems" don't directly relate to the formation of beads of water in this context.

(C) "When the water is too cold" might make sense for freezing, but we're dealing with condensation, not freezing.

(D) "Water spills" just refers to water being spilled, not the condensation process.

(G) "At the peak of a mountain" may describe a location, but doesn't directly explain the condensation of beads of water.

(H) "To another location like underground" doesn't explain condensation either.

Step 5: Focus on plausible answers.

(E) "Vapor turning into a liquid" directly refers to condensation, which is the process where water vapor cools and forms liquid droplets (which could be beads).

(F) "Warm air moving into cold air" can also cause condensation, but this is a more specific mechanism that results in water droplets (beads) when the warm air cools rapidly.

Final answer: [[E]]

Step 1: [[correct]]

Step 2: [[correct]]

Step 3: [[correct]]

Step 4: [[correct]]

Step 5: [[correct]]

Judge: [[]]

Table 16: Critic Direct Example.

Prompt: Please evaluate the reasoning process below carefully, and analyze step by step whether any errors occur in the steps. If all the steps are correct, please output [[]]. You may refer to the following example.

Question: What type of water formation is formed by clouds?

Answer Choice: (A) pearls (B) streams (C) shells (D) diamonds (E) rain (F) beads (G) cooled (H) liquid

Step 1: Understanding the question. The question asks, "What type of water formation is formed by clouds?" We need to determine the correct type of water formation typically associated with clouds.

Step 2: Analyzing the options. The options provided are: (A) pearls (B) streams (C) shells (D) diamonds (E) rain (F) beads (G) cooled (H) liquid

Step 3: Identify the key concept. Clouds are made up of tiny water droplets or ice crystals that have condensed from vapor in the atmosphere. These droplets can eventually combine and fall to the ground as precipitation.

Step 4: Eliminate irrelevant options. Pearls, streams, shells, diamonds, beads, and cooled do not seem to be directly related to water formation associated with clouds.

Step 5: Focus on the most plausible option. "Rain" is the most likely type of water formation that clouds produce, as clouds release water in the form of rain when the droplets or ice crystals become heavy enough to fall.

Final answer: [[E]]

Step 1: Understanding the question. The question asks about the type of water formation formed by clouds. This step is correct. [[correct]]

Step 2: Analyzing the options. This step is also correct in presenting the list of options. [[correct]]

Step 3: Identify the key concept. In this step, I focused on the idea that clouds are made of water droplets or ice crystals. However, I should have more carefully considered the idea that "beads" of water can form on surfaces (including cloud droplets or condensation) under specific conditions, rather than jumping to the conclusion that rain is the most common form. This step is correct. [[correct]]

Step 4: Eliminate irrelevant options. Here's where the mistake lies: I eliminated "beads" prematurely. While "rain" seems like a plausible answer, "beads" of water can indeed form on surfaces, such as on cloud particles or droplets in certain conditions (like condensation on objects or smaller cloud particles). This step is wrong. [[wrong]]

Step 5: Focus on the most plausible option. It dismissed "beads" too quickly. "Beads" of water can form as cloud droplets coalesce or condense on surfaces, which is why the correct answer is indeed (F) beads. This step is correct. [[correct]]

Judge: [[4]]

Question: Where do beads of water come from?

Answer Choice: (A) Too much water (B) underground systems (C) When the water is too cold (D) Water spills (E) Vapor turning into a liquid (F) Warm air moving into cold air (G) At the peak of a mountain (H) To another location like underground

Step 1: Understanding the question.

The question asks where beads of water come from. Beads of water are typically small droplets that form when water vapor condenses.

Step 2: Analyzing the options.

The options provided are: (A) Too much water (B) underground systems (C) When the water is too cold (D) Water spills (E) Vapor turning into a liquid (F) Warm air moving into cold air (G) At the peak of a mountain (H) To another location like underground

Step 3: Focus on key concepts.

Beads of water typically form when water vapor cools and condenses into liquid droplets.

Step 4: Eliminate irrelevant options.

(A) "Too much water" is not a correct explanation for beads of water. (B) "Underground systems" don't directly relate to the formation of beads of water in this context. (C) "When the water is too cold" might make sense for freezing, but we're dealing with condensation, not freezing. (D) "Water spills" just refers to water being spilled, not the condensation process. (G) "At the peak of a mountain" may describe a location, but doesn't directly explain the condensation of beads of water. (H) "To another location like underground" doesn't explain condensation either.

Step 5: Focus on plausible answers.

(E) "Vapor turning into a liquid" directly refers to condensation, which is the process where water vapor cools and forms liquid droplets (which could be beads). (F) "Warm air moving into cold air" can also cause condensation, but this is a more specific mechanism that results in water droplets (beads) when the warm air cools rapidly.

Final answer: [[E]]

Step 1: Understanding the question. The question is clear, and there's no error here. This step is correct. [[correct]]

Step 2: Analyzing the options. This step is also correct. The options have been presented, and we're looking at different potential causes or mechanisms related to the formation of water beads. This step is correct. [[correct]]

Step 3: Focus on key concepts. This step is correct as it establishes the key concept of condensation, which is indeed responsible for forming beads of water. This step is correct. [[correct]]

Step 4: Eliminate irrelevant options. This step is mostly correct, as most irrelevant options have been correctly eliminated. However, (C) might have been a bit misleading. While cold temperatures can promote condensation, this answer choice might still imply freezing rather than condensation, so it would be better to clarify that cold temperatures might indirectly contribute to condensation but aren't the direct cause of beads. This step is correct. [[correct]]

Step 5: Focus on plausible answers. This step is mostly correct. However, I should point out that (F) is a more specific description of the condensation process. It's technically a valid mechanism for water beads to form. So, (F) is not incorrect, but (E) is a more general and direct answer. This step is correct. [[correct]]

Judge: [[]]

Table 17: Critic CoT Example.

Prompt:

- Please evaluate the reasoning step below carefully. You may refer to the previous steps, chosen reward criteria, and the provided scores. For each reasoning step, you must select the appropriate process reward criteria and provide a detailed explanation of the reasoning behind your decision.
- If the selected reward criteria are too coarse-grained or not fully applicable, you should generate more specific fine-grained process reward criteria and justify your choice.
- Each reasoning step must be scored on a scale from 1 to 10 based on the chosen criteria. 1 represents a very poor response, and 10 represents an excellent response. Be sure to clearly explain how you arrived at the score.
- Finally, output the evaluation score in the format: Score: [[X]], where X is the score for the corresponding criterion.

Reward Tree:

- **Coarse-grained Process Criteria:** *Logical Consistency, Mathematical Manipulation, Conceptual Understanding, . . .*
- **Fine-grained Process Reward:** *Correct Application of Formulae, Physical Interpretation, Application of Kepler's Third Law, Energy Interpretation, Relevance to Physical Principles, . . .*

Problem: Consider a satellite that is in a circular orbit around the Earth. The speed of the satellite decreases. What will happen to the satellite's orbital radius and period of revolution? Please justify your answer using principles of physics.

- **Step 1:** The gravitational force provides the centripetal force for circular motion:

$$F_g = \frac{GMm}{r^2} = \frac{mv^2}{r},$$

where M is Earth's mass, m is the satellite's mass, v is orbital speed, and r is orbital radius.

Matching Criteria:

- *Correct Application of Formulae:* The formula $F_g = \frac{GMm}{r^2} = \frac{mv^2}{r}$ is correctly written and applied.
- *Physical Interpretation:* The balance between the gravitational force and centripetal force is clearly described, highlighting how these forces keep the satellite in orbit.

Score:

- *Correct Application of Formulae:* [[10]]
- *Physical Interpretation:* [[9]]

- **Step 2-4:** . . .
- **Step 5:** Applying the escape velocity formula gives

$$v = \sqrt{\frac{2GM}{r}},$$

the new speed corresponds to $r_{\text{new}} = \frac{2GM}{v^2}$.

Matching Criteria:

- *Conceptual Misunderstanding:* The step does not recognize that escape velocity does not directly affect the orbital speed in the problem, leading to the use of an irrelevant formula.
- *Correct Application of Formulae:* The escape velocity formula is incorrectly applied, as it is not relevant to the orbital dynamics in this context.

Score:

- *Conceptual Misunderstanding:* [[2]]
- *Correct Application of Formulae:* [[1]]

-
- **Step 6:** The total mechanical energy

$$E = -\frac{GMm}{2r}$$

increases as v decreases, so r decreases to conserve energy.

Table 18: An example illustrating the DG-PRM process.

Criteria: *Logical Consistency**How to Evaluate:*

- Does the explanation logically progress from one step to the next? Are the connections between orbital speed, radius, and period of revolution clearly explained?
- Does the explanation follow a coherent line of reasoning that reflects an understanding of how physical systems behave?

Questions to Ask:

- Does each step follow logically from the one before? Is there a clear chain of reasoning connecting the steps?

Scoring (1-10):

- 1-3: The explanation is internally inconsistent, with major logical flaws or contradictions.
- 4-6: There are some logical inconsistencies, but the overall reasoning is somewhat coherent.
- 7-8: The explanation is mostly logically consistent, with only minor lapses.
- 9-10: The reasoning is sound, consistent, and flows logically from one point to the next.

Criteria: *Grammar and Writing Style**How to Evaluate:*

- Is the writing free of grammatical errors, spelling mistakes, or awkward sentence structures?
- Does the response maintain a professional and academic tone? Is it easy to read and well-structured?

Questions to Ask:

- Are there any noticeable grammar or spelling issues that detract from the clarity of the explanation?
- Is the sentence structure varied and easy to follow, without excessive repetition or awkward phrasing?

Scoring (1-10):

- 1-3: Numerous grammar or spelling errors, making the explanation difficult to read. The tone may be informal or inappropriate.
- 4-6: Some grammar or spelling errors, but the explanation is still understandable. The tone may occasionally feel inconsistent or too informal.
- 7-8: Few grammar or spelling errors, and the tone is generally appropriate for an academic setting. The writing is mostly clear.
- 9-10: The writing is grammatically correct, well-structured, and professional in tone, with no issues that impede understanding.

Criteria: *Clarity and Precision**How to Evaluate:*

- Is the explanation clear and easy to understand? Are there any overly complex statements, jargon, or vague explanations?
- Are key terms defined and explained appropriately? Is the reasoning direct and to the point?

Questions to Ask:

- Is the explanation concise yet complete, avoiding unnecessary complexity?
- Are key concepts explained in clear and simple terms, avoiding confusion or ambiguity?

Scoring (1-10):

- 1-3: The explanation is difficult to follow, with vague or confusing language.
- 4-6: The explanation is somewhat clear but may be too wordy or not explained in an easily understandable manner.
- 7-8: The explanation is mostly clear and concise, with a few areas of ambiguity.
- 9-10: The explanation is highly clear, well-organized, and precise.

Table 19: Exemplars of coarse-grained process reward criteria.

Criteria: *Relevance to Physical Principles*

How to Evaluate:

- Does the explanation rely on correct and relevant physical principles (e.g., conservation of energy, gravitational forces, orbital mechanics)?
- Are the concepts of orbital speed, radius, and gravitational force properly connected and discussed?

Questions to Ask:

- Does the explanation consider the principles of orbital dynamics and their interrelationships (gravitational force, centripetal force, energy conservation)?
- Are the satellite's movements explained in terms of valid physical laws (such as Newton's laws or Kepler's laws)?

Scoring (1-10):

- 1-3: The explanation is based on incorrect or irrelevant principles, and physics concepts are misapplied or entirely omitted.
- 4-6: Some relevant physical principles are mentioned, but the overall understanding is incomplete or partially incorrect.
- 7-8: Most physical principles are correctly applied, though there may be minor gaps or inaccuracies.
- 9-10: The explanation is rooted in solid, accurate physical principles with correct applications throughout.

Criteria: *Correct Application of Formulae*

How to Evaluate:

- Are the correct equations applied, such as the gravitational force equation, the relation between orbital radius and speed, and Kepler's laws?
- Are any rearrangements of formulas mathematically correct and used appropriately?

Questions to Ask:

- Are the correct formulas used to describe the satellite's motion and interactions with Earth?
- Are these formulas manipulated correctly?

Scoring (1-10):

- 1-3: Incorrect or missing application of key formulas, major mathematical errors.
- 4-6: Some formulas are applied correctly, but there are minor errors or inconsistencies in application or rearrangement.
- 7-8: Most formulas are applied correctly with few minor errors.
- 9-10: All formulas are applied accurately with correct mathematical manipulations.

Criteria: *Identification and Correction of Missteps*

How to Evaluate:

- Does the explanation properly identify any errors in reasoning or misapplications of principles?
- Is the correction provided logical, accurate, and addressing the root cause of the error?

Questions to Ask:

- Does the evaluation recognize and clearly point out the key mistake(s)?
- Is the correction valid and helpful in addressing the mistake? Does it improve the overall understanding of the problem?

Scoring (1-10):

- 1-3: The error is either not identified or misidentified, and the correction is either incorrect or irrelevant.
 - 4-6: The error is identified, but the correction may be incomplete or only partially helpful.
 - 7-8: The error is clearly identified, and the correction is mostly accurate and helpful.
 - 9-10: The error is accurately identified, and the correction provides a clear, valid, and insightful improvement.
-

Table 20: Exemplars of fine-grained process reward criteria (Part I).

Criteria: *Calculation Accuracy**How to Evaluate:*

- Are any numerical or formulaic calculations done correctly, including the proper handling of constants and units?
- If the problem involves mathematical manipulations, is the process free from errors, such as sign mistakes, incorrect square roots, or units mismatch?

Questions to Ask:

- Are all the required calculations correct, with appropriate rounding or approximations applied?
- Are the final results consistent with expectations for orbital mechanics problems (e.g., orbital speed, radius, or period)?

Scoring (1-10):

- 1-3: Numerous calculation errors or mistakes in handling units, resulting in a fundamentally incorrect solution.
- 4-6: Some calculation errors or mismanagement of units, but the overall approach is understandable.
- 7-8: Calculations are mostly accurate, with only minor errors (if any). Units and dimensional analysis are generally correct.
- 9-10: All calculations are accurate and precise, with no errors. Units and dimensions are handled correctly.

Criteria: *Depth of Explanation**How to Evaluate:*

- Does the explanation delve into the key physical concepts and provide a deeper understanding of the underlying physics?
- Is there a thorough discussion of why certain steps are taken and how they relate to the overall problem?

Questions to Ask:

- Does the response merely state the answer, or does it explain the reasoning behind the steps in a comprehensive way?

Scoring (1-10):

- 1-3: The explanation is superficial and lacks a deeper understanding of the principles involved.
- 4-6: The explanation covers the main ideas but does not delve deeply into the concepts behind them.
- 7-8: The explanation shows a solid understanding and provides a good level of depth in reasoning and concept application.
- 9-10: The explanation is deeply insightful, showing an advanced understanding of the concepts and offering a thorough breakdown of the solution.

Criteria: *Understanding of Edge Cases or Special Scenarios**How to Evaluate:*

- Does the response consider potential edge cases or scenarios that might challenge the solution (e.g., extreme values for speed or radius)?

Questions to Ask:

- Does the student acknowledge any potential anomalies or exceptions to the solution, such as very high or low speeds, or situations involving non-ideal conditions (e.g., non-spherical orbits)?
- Does the explanation provide any commentary on how these factors could affect the results?

Scoring (1-10):

- 1-3: The explanation ignores edge cases or unusual scenarios, and assumes the solution applies universally without qualification.
- 4-6: The student touches on edge cases but only briefly or in an incomplete manner.
- 7-8: The student acknowledges and explains edge cases, offering reasonable insights into their impact on the solution.
- 9-10: The student thoroughly discusses edge cases and special scenarios, providing a nuanced understanding of the problem's limitations.

Table 21: Exemplars of fine-grained process reward criteria (Part II).

Instruction: Please carefully analyze each step to determine whether it is correct or incorrect. If you are unsure, you can consult online resources such as Google or AI assistants to help verify the correctness of the information. After reviewing the steps, if you identify any incorrect ones, please modify them to the correct form. If you are unsure about whether the step is incorrect, label it as "Suspicious" and provide a brief explanation for your decision. Ensure that all corrections are clearly stated.
Below is an example.

Problem: Consider a satellite that is in a circular orbit around the Earth. The speed of the satellite decreases. What will happen to the satellite's orbital radius and period of revolution? Please justify your answer using principles of physics.

- **Step 1:** The gravitational force provides the centripetal force for circular motion:

$$F_g = \frac{GMm}{r^2} = \frac{mv^2}{r},$$

where M is Earth's mass, m is the satellite's mass, v is orbital speed, and r is orbital radius.

Evaluation: *Correct* (The equation is a valid representation of the balance between gravitational force and centripetal force for circular motion.)

- **Step 2:** Rearranging the equation gives

$$v = \sqrt{\frac{GM}{r}},$$

showing orbital speed decreases as r increases.

Evaluation: *Correct* (This is a correct rearrangement of the formula, indicating that as the radius increases, the orbital speed decreases.)

- **Step 3:** If v decreases, the gravitational force exceeds the required centripetal force, causing the satellite to spiral inward to a smaller orbit.

Evaluation: *Incorrect* (If v decreases, the satellite would initially move outward, not spiral inward. The gravitational force does not exceed the centripetal force; instead, the satellite adjusts its orbit according to the conservation of energy.) **Correction:** The satellite's orbital radius increases as its speed decreases, so the orbit expands outward, not inward.

- **Step 4:** Using Kepler's third law, $T^2 \propto r^3$, a larger r implies a longer orbital period T .

Evaluation: *Correct* (This is a correct interpretation of Kepler's third law, indicating that as the orbital radius increases, the orbital period also increases.)

- **Step 5:** Applying the escape velocity formula gives

$$v = \sqrt{\frac{2GM}{r}},$$

the new speed corresponds to $r_{\text{new}} = \frac{2GM}{v^2}$.

Evaluation: *Incorrect* (This step incorrectly applies the escape velocity formula. The escape velocity is not directly related to orbital speed in this context, and the equation for r_{new} is not correct.)

Correction: The escape velocity formula should not be used here. Instead, focus on the relationship between orbital speed and radius.

- **Step 6:** The total mechanical energy

$$E = -\frac{GMm}{2r}$$

increases as v decreases, so r decreases to conserve energy.

Evaluation: *Incorrect* (The mechanical energy of a satellite in orbit is negative, and as the speed decreases, the energy decreases, causing the satellite to move to a higher orbit, not a lower one.)

Correction: As v decreases, the mechanical energy becomes less negative, which leads to an increase in the orbital radius.

Tips:

1. Carefully read and analyze the step. Verify its correctness based on physics principles and equations.
 2. If the step is correct, label it as **Correct**.
 3. If the step is incorrect, label it as **Incorrect**, and provide a brief explanation of why it is wrong. Then modify it to the correct version.
 4. If you are unsure whether the step is incorrect, label it as **Suspicious** and provide a brief explanation for your decision.
 5. Ensure that all modifications are clearly stated and the rationale for the correction is clear.
-

Table 22: Instructions for annotators to modify incorrect steps and provide the correct version.

Instruction: For each step in the problem and solution analysis, please evaluate whether the criterion is well-suited to assess the correctness of the step. Please categorize the criteria as “Good”, “Ordinary”, or “Bad” based on the following descriptions:

Good Criterion:

- The criterion is highly relevant and focuses on evaluating the most important aspects of the response.
- It aligns well with the task and encourages the generation of high-quality, focused outputs.
- The criterion is neither too broad nor too narrow and allows for easy application and optimal model adjustment.
- It is practical, straightforward, and well-defined, ensuring consistent and reproducible evaluations.

Label: [Good]

Ordinary Criterion:

- The criterion is somewhat specific but lacks depth or precision to evaluate all aspects effectively.
- It may miss some key considerations or be vague in parts.
- While it provides guidance, it may not lead to optimal model adjustments due to vagueness or complexity.
- It is somewhat practical but lacks clarity in some areas, leading to inconsistent evaluations.

Label: [Ordinary]

Bad Criterion:

- The criterion is irrelevant or focuses on factors that do not contribute meaningfully to the evaluation.
- It is overly complex, ambiguous, or misleading, making it difficult to apply.
- It encourages wrong kinds of responses or evaluates irrelevant aspects of the output.
- It is impractical, rigid, or overly subjective, leading to arbitrary or inconsistent evaluations.

Label: [Bad]

Problem: Consider a satellite that is in a circular orbit around the Earth. The speed of the satellite decreases. What will happen to the satellite’s orbital radius and period of revolution? Please justify your answer using principles of physics.

- **[Step 1-3]:** . . .
- **Step 4:** Using Kepler’s third law, $T^2 \propto r^3$, a larger r implies a longer orbital period T .

Criteria: *Correct Application of Formulae*

How to Evaluate:

- Are the correct equations applied, such as the gravitational force equation, the relation between orbital radius and speed, and Kepler’s laws?
- Are any rearrangements of formulas mathematically correct and used appropriately?

Questions to Ask:

- Are the correct formulas used to describe the satellite’s motion and interactions with Earth?
- Are these formulas manipulated correctly?

Scoring (1-10):

- 1-3: Incorrect or missing application of key formulas, major mathematical errors.
- 4-6: Some formulas are applied correctly, but there are minor errors or inconsistencies in application or rearrangement.
- 7-8: Most formulas are applied correctly with few minor errors.
- 9-10: All formulas are applied accurately with correct mathematical manipulations.

Label: [Good]

Table 23: Instructions for annotators to assess the accuracy of automated validator filtering and the rationality of DG-PRM reward allocation.