# Beyond Correctness: Harmonizing Process and Outcome Rewards through RL Training

**Chenlu Ye**[1,2*]   **Zhou Yu**[1]   **Ziji Zhang**[1]   **Hao Chen**[1]   **Narayanan Sadagopan**[1]

**Jing Huang**[1]   **Tong Zhang**[2]   **Anurag Beniwal**[1]

[1]Amazon      [2]University of Illinois Urbana-Champaign

## Abstract

Reinforcement learning with verifiable rewards (RLVR) has emerged to be a predominant paradigm for mathematical reasoning tasks, offering stable improvements in reasoning ability. However, Outcome Reward Models (ORMs) in RLVR are too coarse-grained to distinguish flawed reasoning within correct answers or valid reasoning within incorrect answers. This lack of granularity introduces noisy and misleading gradients significantly and hinders further progress in reasoning process quality. While Process Reward Models (PRMs) offer fine-grained guidance for intermediate steps, they frequently suffer from inaccuracies and are susceptible to reward hacking.

To resolve this dilemma, we introduce **PRocess cOnsistency Filter (PROF)**, an effective data process curation method that harmonizes noisy, fine-grained process rewards with accurate, coarse-grained outcome rewards. Rather than naively blending PRM and ORM in the objective function (Zou et al., 2025), PROF leverages their complementary strengths through consistency-driven sample selection. Our approach retains correct responses with higher averaged process values and incorrect responses with lower averaged process values, while maintaining positive/negative training sample balance. Extensive experiments demonstrate that our method not only consistently improves the final accuracy over $4\%$ compared to the blending approaches, but also strengthens the quality of intermediate reasoning steps. Codes and training recipes are available at https://github.com/Chenluye99/PROF.

## 1 Introduction

Verifiable rewards have spurred the widest attention recently because they can reliably improve the performance on reasoning tasks with easily verifiable outcomes, such as mathematical and coding problems (Cobbe et al., 2021; Jaech et al., 2024; Shao et al., 2024; Xiong et al., 2025b). However, since the verifiers can only verify the outcome results, the rewards are too sparse and coarse to measure and supervise the reasoning quality in intermediate steps. For instance, if a correct answer contains flawed logic, Outcome Reward Models (ORMs) cannot distinguish it from a completely correct response. We present a classic example from the training data in Table 1, which has fundamentally invalid reasoning but happens to obtain the correct answer. Incorporating such flawed examples into the training process introduces unreliable gradients, leading to significant instability and misguided learning. Moreover, the quality and interpretability of the Chain of Thought (CoT) are crucial for a model's practical reasoning ability, not just the accuracy of the final answer (Zhu et al., 2025; Lyu et al., 2023; Yeo et al., 2024).

---

*Email: chenluy3@illinois.edu

PROF: Process Consistency Filter

> **Problem:** There are coins of $1, 2, 3$, and $5$ cruzeiros, one of each denomination, which weigh $1, 2, 3$, and $5$ grams respectively. One of them is counterfeit, differing in weight from the normal one (but it is not known whether it is heavier or lighter). How can the counterfeit coin be identified with the minimum number of weighings using a balance scale without weights?
>
> ---
>
> **Step 1: Introduction.**
> We need to identify a counterfeit coin with an unknown weight difference from ...
>
> **Step 2 - 6: The Flawed Weighing Logic and Code.**
> 1. Weigh coins of **1 gram and 2 grams** against coins of **3 grams and 5 grams**. 2. Based on the result of the first weighing: - If the left side is heavier, then the counterfeit coin is either 1 gram or 3 grams. - If the right side is heavier, then the counterfeit coin is either 2 gram or 5 grams. - If both sides are equal, then the counterfeit coin is either 2 gram or 5 grams. 3. For the second weighing, ...
> Let's illustrate this with code and ensure that the steps are correct ... output: (1, 'heavier').
>
> > *Analysis: **Incorrect.** The proposed first weighing of {1g, 2g} (total 3g) against {3g, 5g} (total 8g) is fundamentally **invalid**. A balance scale requires comparing groups of equal nominal weight. Because this weighing is unbalanced, all conclusions drawn from it are baseless, including the Python code.*
>
> **Step 7, 8, 9: Summary and Final Answer.**
> Here is a summary of the steps: 1. Weigh coins of 1 gram and 2 grams against coins of 3 grams and 5 grams...
> 2. For the second weighing, weigh the 1 gram coin against a known genuine coin (3 grams)...
> Thus, the minimum number of weighings required to identify the counterfeit coin is $\boxed{2}$.
>
> > *Analysis: **Correct Answer, Invalid Justification.** The summary repeats the flawed weighing strategy. The final answer is correct but is completely unsupported by the fallacious reasoning presented.*

Table 1: An Example of Reaching a Correct Result Through Flawed Reasoning

This limitation of ORMs can be addressed by using LLM-as-a-judge or Monte-Carlo (MC) estimation to provide step-wise judgments or values (Wang et al., 2023; Zheng et al., 2024). However, the cost of inferring LLM step-wise judgments or MC estimation at each iteration during online training is so high. Hence, it is inefficient and expensive to infer the step-wise scores or values for online training. Alternatively, an efficient solution is to use the pre-trained Process Reward Models (PRMs) (Lightman et al., 2023; Zhang et al., 2025), but these models often suffer from misspecification and distribution shift due to the limitation of offline training data. Especially in boundary cases where the policy encounters difficult problems and produces rarely seen responses, PRMs often fail to judge them correctly, thus leading to severe reward hacking (Michaud et al., 2020; Tien et al., 2022). Even if some works (Zha et al., 2025; Cui et al., 2025) attempt to co-train the policy and PRMs online, they can only train in implicit ways such as using implicit generative reward or aligning process rewards with outcomes.

Although numerous works have made enormous efforts to train PRMs offline or online, the problem of effectively coordinating PRMs with outcome-verifiable rewards remains largely underexplored. Existing approaches typically combine process and outcome rewards in a simple weighted manner (Zha et al., 2025; Cui et al., 2025; Zou et al., 2025), which is vulnerable for reward hacking due to the noises and misspecification in PRMs. Therefore, in this paper, instead of developing another PRM, we focus on how to robustly integrate a pre-trained PRM into the online training process, i.e.,

*How to harmonize the accurate but coarse-grained ORMs with fine-grained but noisy Process Reward Models (PRMs) in Reinforcement Learning (RL)?*
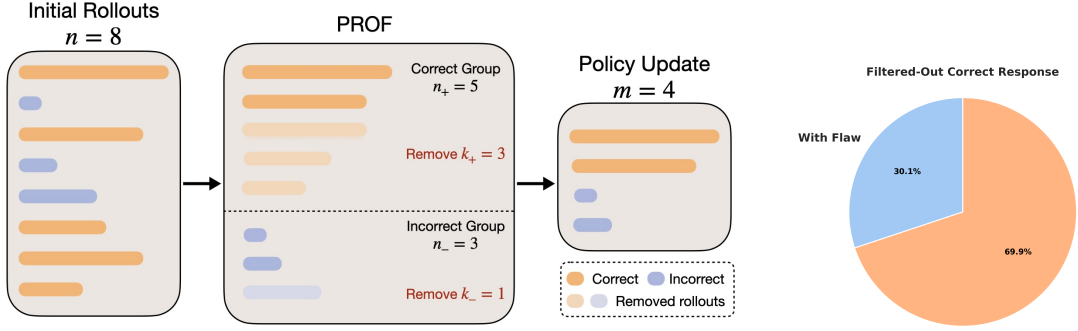
Figure 1: Left: Visualization of PROF Algorithm 1, where the length of each rectangle represents values of process rewards averaged over steps for each rollout. After generating $n$ rollouts and process rewards, PROF ranks the correct and incorrect group separately according to PRM-ORM consistency, so for the correct group, the longer items are kept; for the incorrect group, the shorter items are kept. The number to remove is to balance correct and incorrect ratio. Right: Fraction of flawed-reasoning responses judged by LLM among the filtered-out correct responses.

In this work, instead of fine-tuning another PRM, we answer this question with a **PRocess cOnsistency Filtering (PROF)** framework, a data curation strategy based on process-outcome consistency. PROF over-samples more responses at training time, and then, ranks and filters the responses by the consistency between their PRMs and ORMs. Specifically, it removes samples where the process and outcome signals conflict—such as correct responses derived from flawed reasoning, or incorrect responses that contain sound reasoning steps. By filtering out these inconsistent samples, PROF eliminates conflicting and noisy gradients. Furthermore, observing that correct and incorrect responses have different consistency distributions, we rank each group separately to maintain a balanced training ratio. PROF is a modular framework that can be combined with RL algorithms like Group Relative Policy Optimization (GRPO) for online training.

We conduct extensive experiments to validate the improvement of PROF-GRPO on both outcome accuracy and process reasoning quality at diverse math reasoning benchmarks using both Qwen (Yang et al., 2024) and LLaMA (Dubey et al., 2024) models. To summarize, we highlight our key contributions as follows:

- We propose **PRocess cOnsistency Filtering (PROF)** to robustly integrate noisy Process Reward Models (PRMs) with Outcome Reward Models (ORMs). Compared to the GRPO-type algorithms that only leveraged outcome rewards, our implementation PROF-GRPO can effectively distinguish the inconsistent trajectories, such as correct answers with flawed reasoning steps or incorrect answers with mostly valid steps. Moreover, unlike prior approaches that simply blend PRMs and PRMs, our method only rely on PRMs to rank and filter rather than directly involving them into gradients. This separation essentially avoids reward hacking and entropy collapse, thus achieving stable performance gains throughout training.

- We conduct extensive studies to demonstrate that PROF-GRPO not only increases the final outcome accuracy but also shapes the intermediate reasoning steps and improves the process reasoning quality. Various metrics such as Monte-Carlo estimation, LLM-as-a judge are used to validate that our method enable models to segment reasoning trajectories into **detailed and easy-to-verify** steps.

- We conduct a series of ablation studies to illustrate the importance of separating the correct and incorrect responses during the filtration. Meanwhile, we investigate various ways of calculating the consistency and filtering, and ablate on LLaMA base models for generalization.

## 2  RELATED WORK

**Sample Filtering in Reinforcement Learning for LLM.**   A key challenge in applying reinforcement learning to LLM applications is the imperfection of reward signals. These signals often stem from a learned reward model, such as Reinforcement Learning from Human Feedback (RLHF), or are sparse, delivered only at the end of a trajectory (e.g. RLVR). In Reinforcement Learning from Human Feedback (RLHF), the reward model is trained on human-annotated pairwise comparisons, typically using a Bradley-Terry model (Bradley & Terry, 1952). Due to inherent human disagreement and finite training data, the model develops shortcuts that RL algorithms can exploit (Lin et al., 2023; Eisenstein et al., 2023) to chase for a fake high reward. Consequently, these rewards may not fully align with the underlying intended goals, leading to reward hacking.

Data filtering, a data curation technique, has proven effective in mitigating this issue across various LLM applications with RL. A prominent line of work proposes filtering training pairs based on the reward gap, i.e., the difference in predicted rewards for the chosen and rejected responses (Yuan et al., 2024; Dong et al., 2024; Xiong et al., 2024a; Zhang et al., 2024). The high-level intuition is that a larger reward gap indicates higher model confidence, making these pairs less noisy and more reliable for training when the reward model is well-calibrated. Moreover, Kim et al. (2024); Yu et al. (2025a) further rank and filter the samples by combining their rewards and responses length during the preference learning process.

Similarly, in reasoning tasks employing RL with verifier feedback (RLVR), where rewards are sparse and only provided for the final outcome, filtering is also proven to be helpful. For instance, the simple rejection sampling fine-tuning (Dong et al., 2023; Chen et al., 2025) which discards all incorrect trajectories, serves as a powerful baseline, often approaching the performance of more complex algorithms like GRPO (Dong et al., 2023; Chen et al., 2025; Xiong et al., 2025a). Other methods like (Yang et al., 2024) filters the prompts by difficulty using metrics like pass@k prior to the RL training. Yu et al. (2025b) removes prompts that yield zero gradients during training and dynamically re-generates samples. This technique is known as dynamic sampling and has been rather widely accepted. Xiong et al. (2025a) demonstrates that prompts where all generated responses are incorrect can significantly hurt the performance of the vanilla Reinforce algorithm. They propose an online data filtering strategy based on the correctness reward, showing that a modified Reinforce with filtering (Reinforce-rej) can match or exceed GRPO's performance. Their results suggest that the advantage of GRPO compared to Reinforce is likely due to the implicit data filtering mechanism from the reward shaping. Finally, Xu et al. (2025) proposes to over-sample and keep a subset such that the variance of the rewards in the subset is maximized, which implies that they try to balance the ratio of correct and incorrect responses for reasoning tasks.

In contrast to these methods, which primarily rely on coarse, outcome-based metrics (e.g., final answer correctness, trajectory-level rewards), our approach introduces a more fine-grained filtering mechanism. We leverage process-supervised reward models (PRMs) (Lightman et al., 2023) to evaluate and filter based on the quality of intermediate reasoning steps, and their consistency with ORMs.

**Process-Supervised Reward Models for Fine-Grained Feedback.**   The RLHF focuses on the *trajectory-level comparison* under the Bradley-Terry model. For reasoning-related task, Yang et al. (2024) uses the correctness of the final answer to construct the preference pairs and trains Bradley-Terry reward models for mathematical reasoning. A more widely used approach, termed Outcome Reward Models (ORMs) trains a classifier to predict whether the final answer is correct or not based on the reasoning history. However, Lightman et al. (2023) have shown that Process-Supervised Reward Models (PRMs), which evaluate each intermediate step of a reasoning chain, significantly outperform ORMs, especially for data selection tasks like best-of-n sampling (Lightman et al., 2023). But their approach requires human annotators to label each intermediate steps of the reasoning. Wang et al. (2023) proposes to use Monte-Carlo estimation of the Q value to automatically decide the label. After this, a long line of works proposes to improve the PRMs by

---

**Algorithm 1** Process Consistency Filter (PROF)

---

1: **Input:** Number of rollouts $n$, policy update size $m$, rollout $\{a_1, \ldots, a_n\}$, outcome rewards $\{r_{o,1}, \ldots, r_{o,n}\}$, step number regularization parameter $\lambda, H_\lambda > 0$.

2: Obtain process rewards for each rollout $a_i$ with $H_i$ steps: $(r_i^1, \ldots, r_i^{H_i})$ and compute the trajectory-wise consistency score

$$r_i^{\text{pro}} = \Big[ \frac{1}{H_i} \sum_{h=1}^{H_i} r_i^h - \lambda I(H_i = 1 \text{ or } H_i \geq H_\lambda) \Big] \cdot r_{o,i}. \tag{1}$$

3: Divide rollouts into correct group $\mathcal{G}_+ = \{a_1^+, \ldots, a_{n_+}^+\}$ with $r_{o,i} = 1$ and incorrect group $\mathcal{G}_- = \{a_1^-, \ldots, a_{n_-}^-\}$ with $r_{o,i} = -1$, where $n_+ + n_- = n$.

4: Compute the number to remove $k_+ \in [n_+], k_- \in [n_-]$ in each group:

$$k_+ = \min(n - m, \lceil (\Delta + n - m)/2 \rceil), \ k_- = n - m - k_+, \tag{2}$$

where $\Delta = n_+ - n_-$.

5: Rank $\mathcal{G}_+$ and $\mathcal{G}_-$ by $r^{\text{pro}}$ separately, and keep the samples

$$\mathcal{K}^+ = \{a_i^+ | \text{rank}(a_i^+) \geq n_+ - k_+\}, \ \mathcal{K}^- = \{a_i^- | \text{rank}(a_i^-) \geq n_- - k_-\}.$$

6: **Output:** The kept trajectories $\mathcal{K}^+ \cup \mathcal{K}^-$ with final kept size $m$.

---

generative reward modeling, advanced training technique like RL, and refined engineering practices (Xiong et al., 2024b; Zhang et al., 2025; Khalifa et al., 2025; Zhao et al., 2025; Xiong et al., 2025c). Our work does not focus on improving PRMs but uses the PRMs to supervise the intermediate steps of CoT trajectories for data filtering. We mainly use the Qwen2.5-Math-PRM-7B from Zhang et al. (2025) as it is trained on the distribution of Qwen model and achieves superior performance on ProcessBench (Zheng et al., 2024).

## 3 METHOD

### 3.1 PRELIMINARIES

An LLM is a policy distribution such that given a prompt $x$, it provides the density $\pi(a|x)$ of generating each response $a$. For mathematical reasoning tasks with a binary verifiable reward, there exists a verifier mapping prompt-response pairs $(x, a)$ to a scalar reward $r_o(x, a) \in \{-1, 1\}$. For each prompt, we can generate a group of responses and their corresponding responses with the verifier $\{(a_i, r_{o,i})\}_{i=1}^G$.

**GRPO.** Group Relative Policy Optimization (GRPO) (Shao et al., 2024) is a policy gradient algorithm that simplifies the Proximal Policy Optimization (PPO) (Schulman et al., 2017) by only computing the advantage based on the outcome rewards in a group. Instead of maintaining and updating another value network, GRPO computes the advantage by standardizing the outcome rewards within a group:

$$A_i = \frac{r(x, a_i) - \text{mean}\big(\{r(x, a_j)\}_{j=1}^n\big)}{\text{std}\big(\{r(x, a_j)\}_{j=1}^n\big) + \delta}, \quad i = 1, \ldots, n,$$

where $r(x, a_i)$ is the reward for a given response and $\delta > 0$ is a small constant for numerical stability. This advantage is then incorporated into a clipped surrogate objective function, which is optimized to update the policy from $\pi_{\theta_{\text{old}}}$ to $\pi_\theta$:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}} \left[ \frac{1}{n} \sum_{i=1}^n \frac{1}{|a_i|} \sum_{t=1}^{|a_i|} \min \left( \frac{\pi_\theta(a_{i,<t}|x)}{\pi_{\theta_{\text{old}}}(a_{i,<t}|x)} A_i, \text{clip} \left( \frac{\pi_\theta(a_{i,<t}|x)}{\pi_{\theta_{\text{old}}}(a_{i,<t}|x)}, 1 - \epsilon, 1 + \epsilon \right) A_i \right) \right],$$

where $a_t$ denotes the $t$-th token of response $a$ and $a_{<t}$ denotes $(a_1, \ldots, a_{t-1})$.

Although this approach stabilizes the online policy optimization and is efficient, the sparse reward signal limits further improvement on the intermediate reasoning steps.

**Process Reward Model (PRM).** For a response $a$ containing CoT reasoning, it can be decomposed into multiple reasoning steps $a = (a^1, \ldots, a^H)$. We follow previous works (Zheng et al., 2024; Zhang et al., 2025; Zou et al., 2025) to use a newline as a sign for a new step. For each step $a^h$, the PRM $r^h$ maps it, the previous steps and the prompt $(x, a^{\leq h})$ to a scalar $r^h(x, a^{\leq h})$, where we use the short-hand notation $a^{\leq h} = (a^1, \ldots, a^h)$.

## 3.2 PROF: PROCESS CONSISTENCY FILTER FRAMEWORK

**Algorithm.** We propose our Process Consistency Filter (PROF) in Algorithm 1 to incorporate the consistency of PRMs and ORMs robustly after the rollout phase, and also present a visualization in Figure 1. First, we generate $G$ samples and access the outcome reward. Then, we call the PRM to generate step-wise rewards for each rollout and compute the trajectory-wise consistency score $r^{\mathrm{pro}}$ by taking the mean over the step-wise rewards and adding a step length regularization in equation 1, where $\lambda$ is the regularization parameter and $H_\lambda$ is the threshold for the penalized step number. This regularization is to ensure that samples with no step segments or over-long steps are discarded in the correct group. The samples are divided into two subgroups: $\mathcal{G}_+$ contains the correct samples with $r_o = 1$, and $\mathcal{G}_-$ contains the incorrect samples with $r_o = -1$. Inspired by (Xu et al., 2025), The numbers to discard in each subgroup $k_+, k_-$ are calculated as equation 2 such that the outcome-reward variance of the final kept samples is maximized, which implies that the ratio of correct and incorrect responses should be balanced. After that, we use $r^{\mathrm{pro}}$ to rank and filter the correct group and randomly filter the incorrect group. Finally, we collect the kept $m$ trajectories and use them for policy update.

**Intuition of Consistency Filtration.** To demonstrate that our algorithm can effectively differentiate the inconsistent trajectories, especially those correct answers with flawed reasoning steps, we conduct some understanding experiments. We prompt Qwen2.5-Math-7B-base (Yang et al., 2024) to generate rollouts for 500 problems randomly selected from the training set (introduced in next section) and implement the filtration exactly as Algorithm 1. Then, the filtered-out correct responses are judged by Claude-3-7-sonnet from Anthropic to verify whether they contain flawed steps. We use the prompt in Zhang et al. (2025) and provide the details in Appendix A. As shown in the right plot of Figure 1, 30.1% responses among the filtered-out correct responses are judged to possess flawed reasoning. This indicates that our methods can efficiently distinguish a number of flawed responses and reach consensus with LLM. Furthermore, with human checking those filtered-out correct responses, there are many responses with invalid or even completely wrong reasoning steps but luckily reaching the correct answer. A typical example is presented in Table 1. However, flawed reasoning processes would be entirely missed by a standard ORM.

## 4 EXPERIMENTS

### 4.1 SETUP

**Dataset and Models.** We focus on mathematical reasoning tasks in this work. For online training, we use the prompt set Numina-Math (Beeching et al., 2024) containing nearly 860k math problems with ground-truth answers ranging from Chinese high school math exercises to US and international mathematics Olympiad competition problems. We choose Qwen2.5-Math-1.5B-base, Qwen2.5-Math-7B-base (Yang et al., 2024) as the training base models. For the PRM, we use Qwen2.5-Math-PRM-7B (Zhang et al., 2025) to generate process rewards.

| Model | Algorithm | Math500 | Minerva Math | Olympiad Bench | AIME24 | AMC23 | Average |
|---|---|---|---|---|---|---|---|
| Qwen2.5-Math-1.5B-base | Base | 39.9 | 11.4 | 19.1 | 3.5 | 23.6 | 19.5 |
| | GRPO | 70.3 | 29.1 | 33.0 | 9.0 | 44.5 | 37.2 |
| | Blend | 67.6 | 27.8 | 31.1 | 7.7 | 42.5 | 35.3 |
| | PROF-GRPO | 73.2 | 30.0 | 36.1 | 9.6 | 49.1 | 39.6 |
| Qwen2.5-Math-7B-base | Base | 42.0 | 12.8 | 19.2 | 12.9 | 30.0 | 23.4 |
| | GRPO | 81.6 | 37.2 | 45.5 | 20.6 | 64.4 | 49.9 |
| | Blend | 81.7 | 36.7 | 45.0 | 15.2 | 58.0 | 47.3 |
| | PROF-GRPO | 83.1 | 39.0 | 47.8 | 17.5 | 70.9 | 51.7 |

Table 2: Performance of different algorithms across five benchmarks including Math500 (Hendrycks et al., 2021), Minerva Math (Lewkowycz et al., 2022), Olympiad Bench (He et al., 2024), AMC2023 and AIME2024. We denote Blend-PRM-GRPO by Blend for short. We tune all the algorithms to their best performance and refer the readers to Appendix for the detailed parameter setup. The reported accuracy is average@16 under temperature 1.0.

**Implementation Details.** The implementations are based on the verl framework (Sheng et al., 2025), and we follow most of the parameter settings in verl. Detailedly, we apply the AdamW optimizer with learning rate $1 \times 10^{-6}$. We adopt the clip higher trick (Yu et al., 2025b) that clips the sampling ratio $\pi_\theta/\pi_{\text{old}}$ to an asymmetric range $(1 - \epsilon_{\text{low}}, 1 + \epsilon_{\text{high}})$. Specifically, we set $\epsilon_{\text{low}} = 0.2$, $\epsilon_{\text{high}} = 0.28$ for models started from Qwen2.5-Math-1.5B-base and maintain $\epsilon_{\text{high}} = \epsilon_{\text{low}} = 0.2$ for other cases. In each iteration, we sample 1024 prompts, rollout $n = 4$ responses per prompt for GRPO and $n = 8$ responses for PROF-GRPO. Note that the policy update number for all algorithms is $m = 4$. For the regularization of step numbers in Algorithm 1, we take $\lambda = 10$ and $H_\lambda = 30$. More details are provided in Appendix A and scripts are available in the GitHub repository.

**Evaluation.** The models' performance is evaluated on 5 benchmarks: Math500 (Hendrycks et al., 2021), Minerva Math (Lewkowycz et al., 2022), Olympiad Bench (He et al., 2024), AMC2023[1] and AIME2024[2]. We mainly use average@16 for evaluation, i.e., the accuracy is averaged over 16 responses per prompt under temperature 1.0. The models are allowed to generate 4096 tokens.

## 4.2 MAIN RESULTS

We summarize our main results in Table 2, where Blend denotes a commonly used way that mixes the PRM with outcome rewards (Zha et al., 2025; Cui et al., 2025; Zou et al., 2025). Following (Zou et al., 2025), the PRMs are averaged over steps for each response, weighted by a parameter $\beta$, and added to the outcome rewards. We use the parameter $\beta = 0.8$ according to Table 5 of (Zou et al., 2025). Our main findings are as follows.

**PROF-GRPO Outperforms the Baselines.** As shown in Table 2, our proposed method, PROF-GRPO, consistently outperforms GRPO and Blend-PRM-GRPO over various benchmarks. Specifically, for models starting from Qwen2.5-Math-1.5B-base, PROF-GRPO achieves an average accuracy of 39.6%, surpassing the standard GRPO baseline (37.2%) and the Blend-PRM-GRPO method (35.3%). A similar trend is ob-

---

[1]https://huggingface.co/datasets/math-ai/amc23
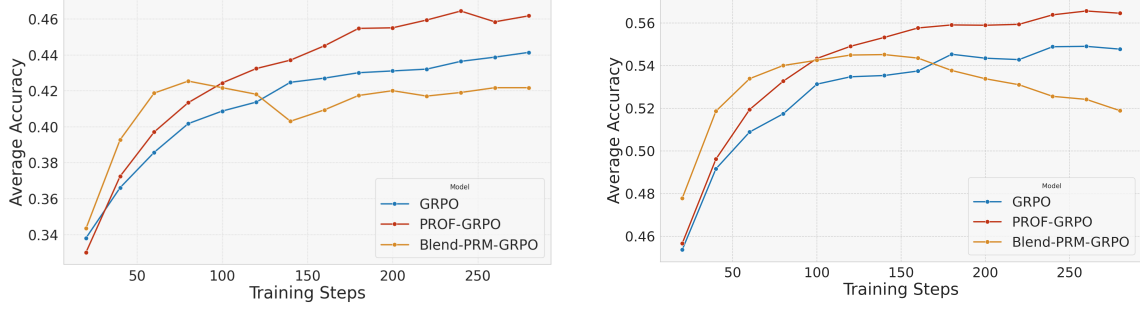[2]https://huggingface.co/datasets/math-ai/aime24

Figure 2: The learning dynamics of PROF-GRPO initialized from Qwen2.5-Math-1.5B-base (left) and Qwen2.5-Math-7B-base (right). The y-axis is the average@16 accuracy and is further averaged on Math500, Minerva Math and Olympiad Bench. For comparison, we also plot the baseline GRPO and directly blending PRM with outcome rewards, denoted as Blend-PRM-GRPO.
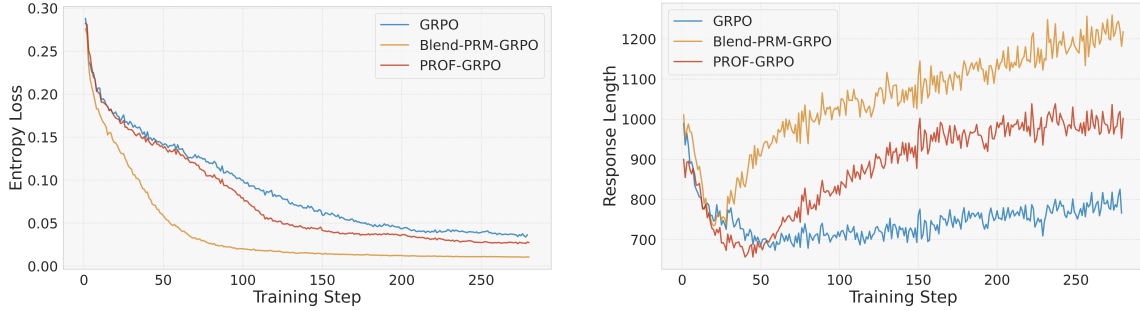


Figure 3: Left: the entropy loss curves of GRPO, Blend-PRM-GRPO and PROF-GRPO; Right: the response length curves of GRPO, Blend-PRM-GRPO and PROF-GRPO, all of which are initialized from Qwen2.5-Math-7B-base.

served with the Qwen2.5-Math-7B-base model, where PROF-GRPO achieves a $51.7\%$ average accuracy, a significant improvement over GRPO's $49.9\%$ and Blend-PRM-GRPO's $47.3\%$. [3]

The learning dynamics in Figure 2 corroborate these findings, illustrating that PROF-GRPO steadily maintains a consistent performance advantage over both GRPO and Blend-PRM-GRPO throughout the training process. Notably, PROF-GRPO achieves faster convergence rate and higher final accuracy than GRPO.

**Filtration Method is Much More Robust than Blending.** We plot the entropy loss and response length curves of GRPO, Blend-PRM-GRPO and PROF-GRPO in Figure 3. Combining it and Figure 2, we observe that Blending-PRM-GRPO suffers from severe rewarding hacking. Figure 3 shows that its entropy collapse quickly towards zero. Simultaneously, its response length in the right figure uncontrollably increases, indicating that the model has learned to game the PRM by over-generating verbose responses and more repetitive steps to get a higher averaged process reward. Therefore, Blend-PRM-GRPO's testing accuracy even falls below GRPO. In contrast, PROF-GRPO maintains gradual and slightly faster decrease in entropy

---

[3]Although PROF-GRPO underperformed relative to GRPO on AIME24 for Qwen2.5-Math-7B-base, this result should be interpreted with caution. Given the dataset's small size of only 30 samples, the performance difference may not be statistically significant.

loss and controllable response length growth. This illustrates that our filtration method effectively leverages the PRM signal while stay robust to reward hacking. We will carefully analyze and compare the quality of intermediate reasoning steps of our method and baselines.

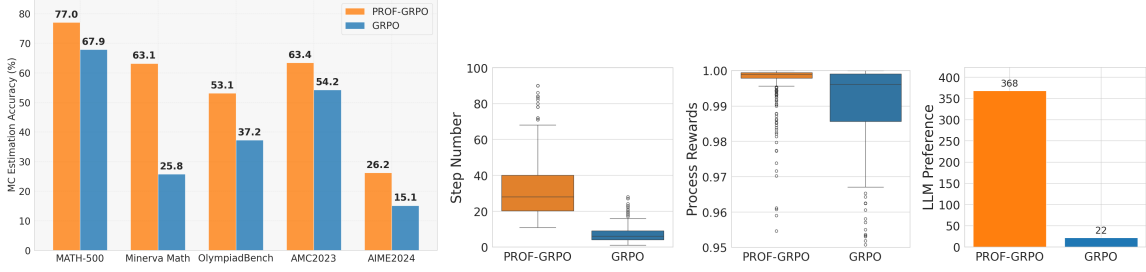## 4.3 HOW PROF SHAPES INTERMEDIATE REASONING STEPS



Figure 4: Reasoning intermediate-steps performance of PROF-GRPO in comparison with GRPO. The most left plot is the Monte Carlo (MC) estimation scores across five benchmarks. The other three are on Math500 under metrics of number of steps (2nd left), the averaged process rewards generated by Qwen2.5-Math-PRM-7B (3rd left), and LLM's preference between two modes' responses (most right).

**Improved Step-wise Value.** To evaluate the quality of intermediate steps, we adopt Monte Carlo (MC) estimation, a commonly used way to estimate the probability of getting to the correct final answer and estimate process values (Wang et al., 2023; Xiong et al., 2024a; Luo et al., 2024). For this analysis, we select problem-response pairs from the test prompts where our method (PROF-GRPO) and GRPO both produced the correct final answer. Both models were initialized from Qwen2.5-Math-7B-base. To estimate the value of each reasoning step, we generate eight independent completions from that point using a temperature of 1.0, and the resulting empirical success rate serves as the MC value.

Our primary finding is that PROF-GRPO achieves a significant improvement in step-wise values compared to GRPO. As shown in the leftmost plot of Figure 4, the average MC estimations across all five benchmarks are consistently higher for our model. The specific improvement gaps are 9.2% on Math500, 37.4% on Minerval Math, 15.9% on Olympiad Bench, 9.2% on AMC2023, and 11.1% on AIME2024, which are much larger than the outcome accuracy gap in Table 2. This finding validates that in addition to improve the outcome accuracy, more importantly, our PROF method substantially improves the quality and consistency of intermediate steps.

**Deeper Analysis on Math500.** We conducted further comparisons on the Math500 benchmark using responses where both models were correct in figure 4. In the second left figure, PROF-GRPO's responses have more reasoning steps. In the third left figure, the PRM used for training assigns higher rewards for PROF-GRPO's responses. In the most right figure, we use Claude-3-7-sonnet from Anthropic to judge which one's reasoning process has more complete and detailed steps, and PROF-GRPO's responses are significantly preferred. The prompt for LLM-as-a-judge is presented in Table A.2. Overall, the key takeaway is that our PROF method reshapes the model's CoT process from complicated and skipping reasoning into **detailed and easy-to-verify** steps. This can be further validated by two examples in Figure 5 and 8.

## 5 ABLATIONS

To understand the superior performance of PROF-GRPO, we conduct a series of ablation studies.

**Problem:** A white dwarf star has an effective temperature, $T_e = 50,000$ degrees Kelvin, but its radius $R_{WD}$ is comparable to that of the Earth. Take $R_{WD} = 10^4$km($10^7$m or $10^9$ cm). Compute the luminosity (power output) of the white dwarf. Treat the white dwarf as a blackbody radiator. Give your answer in units of ergs per second, to two significant figures.

**PROF-GRPO Solution:**

First, let's recall the formula for the luminosity of a blackbody.
The luminosity $L$ of a blackbody is given by the Stefan-Boltzmann law: $L = 4\pi R^2 \sigma T_e^4$
We are given: $T_e = 50,000$K... The Stefan-Boltzmann constant $\sigma$ is approximately $5.67 \times 10^{-8}$W$m^{-2}$K$^{-4}$.
Next, let's substitute the values into the Stefan-Boltzmann law.
$L = 4\pi(10^7)^2 \times 5.67 \times 10^{-8} \times (50000$K$)^4$
Let's calculate the value step by step.
$(10^7)^2 = 10^{14}$, $(50000$K$)^4 = \cdots = 6.25 \times 10^{18}$.
Now, substitute these values back into the equation..., calculate the numerical part
$4\pi \approx 4 \times 3.14159 = 12.56636$,
$12.56636 \times 5.67 \approx 71.148$
$71.148 \times 6.25 \approx 444.675$
$L \approx \approx 4.45 \times 10^{26}$ W
Since 1 W=1 J/s and 1 erg=$10^{-7}$ J, ...
$L \approx 4.45 \times 10^{33}$ erg/s
Therefore, the luminosity of the white dwarf is:
$\boxed{4.5 \times 10^{33} \text{ erg/s}}$ ✅

**GRPO Solution:**

To compute the luminosity of the white dwarf, we can use the Stefan-Boltzmann law
...
The formula for the luminosity $L$ of a blackbody is given by: $L = 4\pi R^2 \sigma T_e^4$,
where $\sigma = 5.67 \times 10^{-8}$W$m^{-2}$K$^{-4}$.
Given: $T_e = 5 \times 10^4$K, $R_{WD} = 10^9$cm.
First, we need to calculate $T_e^4 = (5 \times 10^4)^4 = 625 \times 10^{16} = 6.25 \times 10^{18}$K$^4$
Next, we calculate the surface area of the white dwarf:
$4\pi R_{WD}^2 = 4\pi(10^9$cm$)^2 = 4\pi \times 10^{18}$cm$^2$.
Now we can compute the luminosity:
$L = 4\pi \times 10^{18}$cm$^2 \times 5.67 \times 10^{-5}$ergs$^{-1}$cm$^{-2}$K$^{-4} \times 6.25 \times 10^{18}$K$^4$.
$L = 4\pi \times 5.67 \times 6.25 \times 10^{18} \times 10^{-5} \times 10^{18}$ergs$^{-1}$
$= \cdots = 4.42 \times 10^{33}$ergs$^{-1}$s$^{-1}$ ❌
Rounded to two significant figures, the luminosity of the white dwarf is
$\boxed{4.4 \times 10^{33} \text{ergs}^{-1}}$.

**Blend-PRM-GRPO Solution:**

The luminosity $L$ of a blackbody can be calculated using the Stefan-Boltzmann law:
$L = 4\pi R^2 \sigma T_e^4$, where $\sigma$ is approximately $5.67 \times 10^{-5}$ ergs per square centimeter per second per Kelvin to the fourth power.
$R_{WD} = 10^9$cm.
First, we need to convert the radius from kilometers to centimeters.
First, we calculate $R_{WD}^2 = (10^9$cm$)^2 = 10^{18}$cm$^2$
Next, we calculate...$T_e^4 = 6.25 \times 10^{18}$K$^4$
Now we can substitute these values into the Stefan-Boltzmann law:
$L = 4\pi \times 10^{18}\text{cm}^2 \times 5.67 \times 10^{-5}\text{ergs/cm}^2\text{/s/K}^4 \times 6.25 \times 10^{18}$K$^4$
... We can separate the constants and the powers of 10:
$4\pi \times 5.67 \times 6.25 \times 10^{-5+18}$
...
So, $444.51875 \times 10^{13}$ ❌
Therefore, $L \approx 4.45 \times 10^{15}$ ergs/s
So the luminosity of the white dwarf is
$\boxed{4.45 \times 10^{15} \text{ ergs/s}}$.

Figure 5: A Minerval-Math example to compare distinct intermediate reasoning patterns of PROF-GRPO, vanilla GRPO and Blend-PRM-GRPO. PROF-GRPO presents concrete and correct deduction steps. GRPO's solution skips detailed deduction steps and there are flaws in the calculation precision and the final rounding. Blend-PRM-GRPO has long-winded steps and makes a big mistake in calculating the power.

## 5.1 SEPARATION OF CORRECT AND INCORRECT GROUP

To justify the necessity of handling correct and incorrect samples separately, we conduct an ablation experiment **PROF-GRPO w/o separation**, where the $n$ rollouts are ranked and filtered without separating correct and incorrect responses. To mitigate bias in PRM, each step's PRM is subtracted by the averaged PRM of the batch. Even after centering, the lower left plot in Figure 6 shows that PROF-GRPO w/o Separation has over $2\%$ gap between the training reward after and before the filtration. This indicates that a disproportionate number of negative samples are removed, thereby distorting the ratio of correct to incorrect samples. One explanation is that incorrect responses often contain several correct intermediate steps, thus increasing the averaged PRM over steps and leading to lower consistency. Consequently, incorrect responses exhibit lower consistency than correct ones, especially as the policy model improves over training. In contrast, PROF-GRPO successfully balances the bias by separating the correct and incorrect groups and maintains nearly the same reward before and after filtering.

To further disentangle the contributions of filtering correct versus incorrect samples, we design the following variants of PROF:

(1) Filter-Correct: use PRM consistency to filter the correct group and randomly filter the incorrect group;

(2) Filter-Incorrect: use PRM consistency to filter the incorrect group and randomly filter the correct group;

(3) Filter-Random: only maximize the outcome reward variance and randomly filter both correct and incorrect samples, following Xu et al. (2025).

As shown in the first row of Figure 6, Filter-Correct and PROF-GRPO (Filter-both) achieve comparably best performances among the variants across the 1.5B and 7B models. While Filter-both converges more effi-
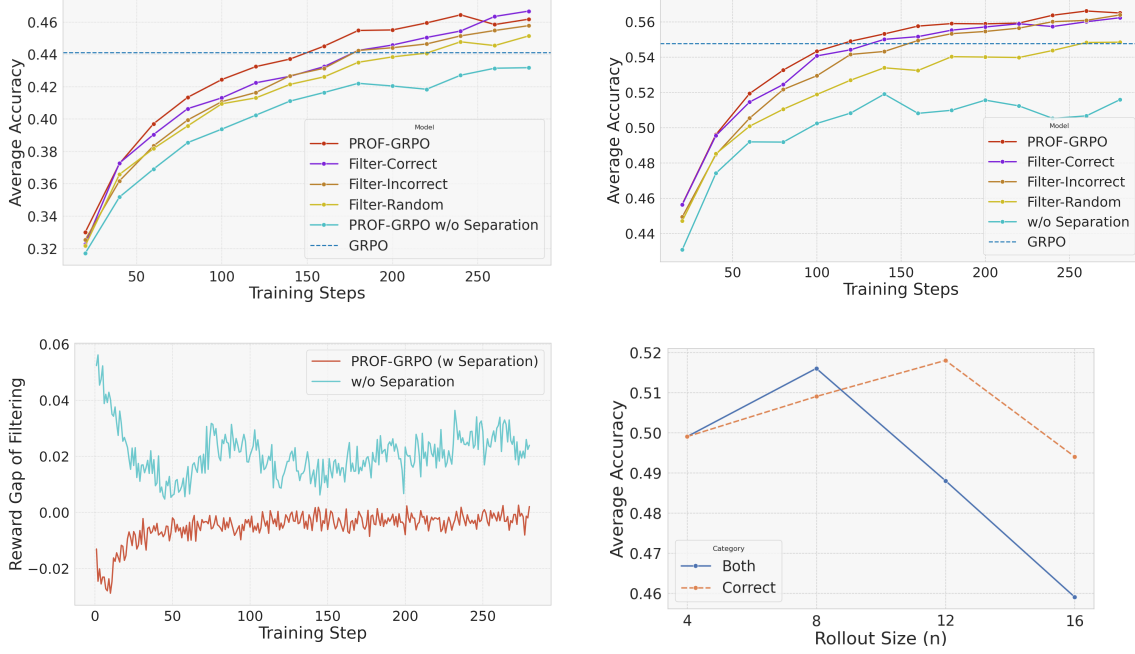
Figure 6: First row: the averaged accuracy over three benchmarks Math500, Minerva Math and Olympiad Bench for PROF-GRPO (Filter both correct and incorrect with PRM consistency) and its variants initialized from wen2.5-Math-1.5B-base (left) and Qwen2.5-Math-7B-base (right). Lower left: the gap between the training rewards after and before the filtering for PROF-GRPO in comparison with not separating correct and incorrect groups (w/o separation). Lower right: the averaged accuracy across all five benchmarks over rollout sizes $n = 4, 8, 12, 16$ for filtering both correct and incorrect groups with PRM consistency (Both) and only the correct group with PRM consistency (Correct).

ciently because it leverages the consistency filtration for both correct and incorrect groups. Filter-incorrect is less efficient and has slightly poorer performance. In contrast, Filter-Random only performs slightly better than GRPO, and w/o Separation performs the worst.

These findings lead to a clear conclusion: separating the correct and incorrect groups is essential to prevent the over-removal of valuable incorrect samples during training. While both Filter-both and Filter-Correct are top-performing strategies, with the former being more efficient, the trade-offs between them will be discussed in the following section. Furthermore, the comparable performance of Filter-both and Filter-Correct indicates that the process quality for correct samples is more crucial than the consistency for incorrect samples during the training process.

## 5.2 EFFECT OF ROLLOUT NUMBERS

We study the scale of rollout numbers $n$ with fixed policy-update number $m = 4$ by varying $n = 4, 8, 12, 16$. The lower-right plot in Figure 6 presents the test accuracy averaged over all five benchmarks for PROF-GRPO (Both) and Filter-Correct (Correct) started from Qwen2.5-Math-7B-base. We observe the performance first increases then decreases as $n$ increases, revealing a trade-off between enhancing process reasoning quality and avoiding reward hacking. Notably, Filter-Correct decreases later (after $n = 12$) because it

only leverages the influence of PRM only in the correct group, indicating that Filter-Correct is more robust when the PRM's influence is higher, like when increasing the scale of ranking and filtering.

## 5.3 VARIANTS OF FILTRATION METHODS

| Algorithm | Math500 | Minerva Math | Olympiad Bench | AIME24 | AMC23 | Average |
|---|---|---|---|---|---|---|
| Mean | **83.1** | **39.0** | **47.8** | 17.5 | **70.9** | **51.7** |
| Minimum | 82.9 | 38.3 | 46.7 | 20.8 | 65.9 | 50.9 |
| Sum | 82.4 | 38.1 | 47.4 | 17.7 | 67.5 | 50.6 |
| Ratio | 81.4 | 36.6 | 45.0 | **24.8** | 65.2 | 50.6 |

Table 3: Performance of different filtration ways in PROF starting from Qwen2.5-Math-7B-base.

In this subsection, we investigate the influence of different computation methods of consistency score $r^{\mathrm{pro}}$ in addition to the mean of PRMs over steps, where Mean denotes averaging over steps in Algorithm 1, Minimum and Sum denotes taking the minimum and sum summation over steps, Ratio denotes filtering while preserving the original positive–negative sample distribution, instead of balancing. As shown in Table 3, the performances of Minimum (50.9%), Sum (50.6%), and Ratio (50.6%) are inferior to the mean. This suggests that the mean provides a more stable estimate of reasoning consistency: unlike the minimum, it is less sensitive to a single poorly scored step, and unlike the summation, it avoids bias towards longer trajectories. Additionally, balancing the correct-incorrect ratio can use data consistency to select the group with more sufficient samples without breaking their balance.

## 5.4 ABLATION STUDY ON BASE MODEL

| Algorithm | Math500 | Minerva Math | Olympiad Bench | AIME24 | AMC23 | Average |
|---|---|---|---|---|---|---|
| Base | 30.0 | 8.8 | 6.1 | 2.3 | 10.6 | 11.6 |
| GRPO | 50.5 | 18.8 | 17.9 | 5.0 | 25.6 | 23.6 |
| PROF-GRPO (Both) | 50.4 | 19.1 | 18.7 | 3.5 | 27.8 | 23.9 |
| PROF-GRPO (Correct) | **52.4** | **19.5** | **19.8** | **6.7** | **28.6** | **25.4** |
| PROF-GRPO (Incorrect) | 49.0 | 18.0 | 17.3 | 5.4 | 23.9 | 22.7 |
| Blend-PRM-GRPO | 37.2 | 13.1 | 9.9 | 1.0 | 17.2 | 15.7 |

Table 4: The test accuracy of different methods initialized from LLaMA-3.2-3B-instruct that is average@16 under temperature 1.0 and further averaged across all the five benchmarks.

To showcase the generalization of our algorithm, we conduct experiments on LLaMA-3.2-3B-instruct (Dubey et al., 2024) that has weaker math-reasoning abilities and more distribution shift since Qwen2.5-Math-PRM-7B is trained on the distribution of Qwen's family. As provided in Table 4, PROF-GRPO with PRM consistency filtering both correct and incorrect groups (Both) achieves 23.9%, marginally outperforming the GRPO baseline (23.6%), while only applying PRM consistency to filter the correct group (Correct) exhibits the strongest (25.4%) performance. Conversely, applying the filter solely to the incorrect group (PROF-GRPO (Incorrect)) is counterproductive, causing accuracy to drop to 22.7%. Blend-PRM-GRPO still scores the worst (15.7%) among all the methods. These results suggest that our PROF methods can consistently outperform baselines across various base models.

For the trade-off between the Both and Correct, we conclude that when the PRM is less reliable or prone to reward hacking (as in this cross-model scenario), the "Correct" method offers more robust improvements by safely constraining the PRM's influence. However, when the PRM is highly reliable and training efficiency is a priority, the "Both" method is recommended.

## 6 CONCLUSION AND FUTURE WORK

In this work, we introduce the Process Consistency Filter (PROF), a novel data curation technique that filters generated responses by the data PRM-ORM consistency, and maintains the balance of correct and incorrect ratios. We demonstrate its effectiveness and robustness in not only consistently improving the accuracy of obtaining correct final answers, but also shaping the policy model to generate more detailed and fine-grained segmented intermediate reasoning steps. Particularly, PROF is a general subsampling and filtration framework without reliance on specific PRMs or the RL algorithms. Thus, the use of Qwen2.5-Math-PRM-7B as the PRM in our experiments is not a limitation. Exploring the integration of PROF with more accurate or diverse PRMs remains an interesting direction for future work. Additionally, how to extend our method to other reasoning tasks, such as coding (Jimenez et al., 2023) and web navigation (Zhou et al., 2023) also deserves to be explored.

## REFERENCES

Edward Beeching, Shengyi Costa Huang, Albert Jiang, Jia Li, Benjamin Lipkin, Zihan Qina, Kashif Rasul, Ziju Shen, Roman Soletskyi, and Lewis Tunstall. Numinamath 7b cot. https://huggingface.co/AI-MO/NuminaMath-7B-CoT, 2024.

Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.

Huayu Chen, Kaiwen Zheng, Qinsheng Zhang, Ganqu Cui, Yin Cui, Haotian Ye, Tsung-Yi Lin, Ming-Yu Liu, Jun Zhu, and Haoxiang Wang. Bridging supervised learning and reinforcement learning in math reasoning. *arXiv preprint arXiv:2505.18116*, 2025.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang, Wendi Li, Bingxiang He, Yuchen Fan, Tianyu Yu, Qixin Xu, Weize Chen, et al. Process reinforcement through implicit rewards. *arXiv preprint arXiv:2502.01456*, 2025.

Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*, 2023.

Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen Sahoo, Caiming Xiong, and Tong Zhang. Rlhf workflow: From reward modeling to online rlhf. *arXiv preprint arXiv:2405.07863*, 2024.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.

Jacob Eisenstein, Chirag Nagpal, Alekh Agarwal, Ahmad Beirami, Alex D'Amour, DJ Dvijotham, Adam Fisch, Katherine Heller, Stephen Pfohl, Deepak Ramachandran, et al. Helping or herding? reward model ensembles mitigate but do not eliminate reward hacking. *arXiv preprint arXiv:2312.09244*, 2023.

Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint arXiv:2402.14008*, 2024.

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.

Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.

Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.

Muhammad Khalifa, Rishabh Agarwal, Lajanugen Logeswaran, Jaekyeom Kim, Hao Peng, Moontae Lee, Honglak Lee, and Lu Wang. Process reward models that think. *arXiv preprint arXiv:2504.16828*, 2025.

Sunnie SY Kim, Q Vera Liao, Mihaela Vorvoreanu, Stephanie Ballard, and Jennifer Wortman Vaughan. " i'm not sure, but...": Examining the impact of large language models' uncertainty expression on user reliance and trust. In *Proceedings of the 2024 ACM conference on fairness, accountability, and transparency*, pp. 822–835, 2024.

Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in neural information processing systems*, 35:3843–3857, 2022.

Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023.

Yong Lin, Hangyu Lin, Wei Xiong, Shizhe Diao, Jianmeng Liu, Jipeng Zhang, Rui Pan, Haoxiang Wang, Wenbin Hu, Hanning Zhang, et al. Mitigating the alignment tax of rlhf. *arXiv preprint arXiv:2309.06256*, 2023.

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. Improve mathematical reasoning in language models by automated process supervision, 2024. URL https://arxiv.org/abs/2406.06592.

Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. Faithful chain-of-thought reasoning. In *The 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (IJCNLP-AACL 2023)*, 2023.

Eric J Michaud, Adam Gleave, and Stuart Russell. Understanding learned reward functions. *arXiv preprint arXiv:2012.05862*, 2020.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. In *Proceedings of the Twentieth European Conference on Computer Systems*, pp. 1279–1297, 2025.

Jeremy Tien, Jerry Zhi-Yang He, Zackory Erickson, Anca D Dragan, and Daniel S Brown. Causal confusion and reward misidentification in preference-based reward learning. *arXiv preprint arXiv:2204.06601*, 2022.

Peiyi Wang, Lei Li, Zhihong Shao, RX Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. *arXiv preprint arXiv:2312.08935*, 2023.

Wei Xiong, Chengshuai Shi, Jiaming Shen, Aviv Rosenberg, Zhen Qin, Daniele Calandriello, Misha Khalman, Rishabh Joshi, Bilal Piot, Mohammad Saleh, et al. Building math agents with multi-turn iterative preference learning. *arXiv preprint arXiv:2409.02392*, 2024a.

Wei Xiong, Hanning Zhang, Nan Jiang, and Tong Zhang. An implementation of generative prm, 2024b.

Wei Xiong, Jiarui Yao, Yuhui Xu, Bo Pang, Lei Wang, Doyen Sahoo, Junnan Li, Nan Jiang, Tong Zhang, Caiming Xiong, et al. A minimalist approach to llm reasoning: from rejection sampling to reinforce. *arXiv preprint arXiv:2504.11343*, 2025a.

Wei Xiong, Hanning Zhang, Chenlu Ye, Lichang Chen, Nan Jiang, and Tong Zhang. Self-rewarding correction for mathematical reasoning. *arXiv preprint arXiv:2502.19613*, 2025b.

Wei Xiong, Wenting Zhao, Weizhe Yuan, Olga Golovneva, Tong Zhang, Jason Weston, and Sainbayar Sukhbaatar. Stepwiser: Stepwise generative judges for wiser reasoning, 2025c. URL https://arxiv.org/abs/2508.19229.

Yixuan Even Xu, Yash Savani, Fei Fang, and Zico Kolter. Not all rollouts are useful: Down-sampling rollouts in llm reinforcement learning. *arXiv preprint arXiv:2504.13818*, 2025.

An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2. 5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.

Wei Jie Yeo, Ranjan Satapathy, Rick Siow Mong Goh, and Erik Cambria. How interpretable are reasoning explanations from prompting large language models? *arXiv preprint arXiv:2402.11863*, 2024.

Ping Yu, Weizhe Yuan, Olga Golovneva, Tianhao Wu, Sainbayar Sukhbaatar, Jason Weston, and Jing Xu. Rip: Better models by survival of the fittest prompts. *arXiv preprint arXiv:2501.18578*, 2025a.

Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025b.

Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*, 3, 2024.

Kaiwen Zha, Zhengqi Gao, Maohao Shen, Zhang-Wei Hong, Duane S Boning, and Dina Katabi. Rl tango: Reinforcing generator and verifier together for language reasoning. *arXiv preprint arXiv:2505.15034*, 2025.

Chuheng Zhang, Wei Shen, Li Zhao, Xuyun Zhang, Lianyong Qi, Wanchun Dou, and Jiang Bian. Policy filtration in rlhf to fine-tune llm for code generation. 2024.

Zhenru Zhang, Chujie Zheng, Yangzhen Wu, Beichen Zhang, Runji Lin, Bowen Yu, Dayiheng Liu, Jingren Zhou, and Junyang Lin. The lessons of developing process reward models in mathematical reasoning. *arXiv preprint arXiv:2501.07301*, 2025.

Jian Zhao, Runze Liu, Kaiyan Zhang, Zhimu Zhou, Junqi Gao, Dong Li, Jiafei Lyu, Zhouyi Qian, Biqing Qi, Xiu Li, and Bowen Zhou. Genprm: Scaling test-time compute of process reward models via generative reasoning, 2025. URL https://arxiv.org/abs/2504.00891.

Chujie Zheng, Zhenru Zhang, Beichen Zhang, Runji Lin, Keming Lu, Bowen Yu, Dayiheng Liu, Jingren Zhou, and Junyang Lin. Processbench: Identifying process errors in mathematical reasoning. *arXiv preprint arXiv:2412.06559*, 2024.

Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.

Dawei Zhu, Xiyu Wei, Guangxiang Zhao, Wenhao Wu, Haosheng Zou, Junfeng Ran, Xun Wang, Lin Sun, Xiangzheng Zhang, and Sujian Li. Chain-of-thought matters: improving long-context language models with reasoning path supervision. *arXiv preprint arXiv:2502.20790*, 2025.

Jiaru Zou, Ling Yang, Jingwen Gu, Jiahao Qiu, Ke Shen, Jingrui He, and Mengdi Wang. Reasonflux-prm: Trajectory-aware prms for long chain-of-thought reasoning in llms. *arXiv preprint arXiv:2506.18896*, 2025.

# A    ADDITIONAL EXPERIMENTAL DETAILS AND RESULTS

## A.1    MAIN EXPERIMENTS

For the rollout stage, we use a temperature of $1.0$ and a top-p value of $1.0$. We set the KL loss coefficient to $0.001$ and entropy loss coefficient to $0.001$. All the models are trained with 8 H100 GPUs. We set the training mini-batch size as 256 and allow the models to generate 4096 tokens per prompt.

## A.2    PROMPT TEMPLATE

We present the template used for LLM to compare step-level reasoning.

---

**Prompt for Finding Reasoning Flaws in Correct Response via LLM-as-a-judge**

Here is the problem and the assistant's solution, which has been broken down into {step} steps.

**Problem:**

**Assistant's Solution:**

Your task is to review each step of the solution in sequence, analyzing, verifying, and critiquing the reasoning in detail. You need to provide the analyses and the conclusion in the following format:
<step>Step 1 Analysis</step>
<step>Step 2 Analysis</step>
... [CONTINUE FOR ALL step steps in the Assistant's Solution] ...
<conclusion>Correct/Incorrect</conclusion>

- When you analyze each step, you should use proper verification, recalculation, or reflection to indicate whether it is logically and mathematically valid. Please elaborate on the analysis process carefully.

- If an error is detected in any step, you should describe the nature and cause of the error in detail, and suggest how to correct the error or the correct approach. Once a step is found to contain any error, stop further analysis of subsequent steps (as they may depend on the identified error) and directly provide the conclusion of "Incorrect."

For instance, given a solution of five steps, if an error or flaw is found in the third step, you should reply in the following format:
<step>Step 1 Analysis</step>
<step>Step 2 Analysis</step>
<step>Step 3 Analysis; since an error or flaw is found here, also provide detailed critique and correction guideline)</step>
<conclusion>Incorrect</conclusion>
Respond with your analyses and conclusion directly.

---

---

**Prompt for Responses Comparison via LLM-as-a-judge**

**System** You are a meticulous, comparison engine. Your ONLY function is to compare the intermediate reasoning steps of the two responses provided to you.
**User** Here is the problem and assistants' two solutions, which have been chunked into steps. You MUST provide preference over the two solutions.
Problem: <prompt>
Assistant's Solution 1: <solution1>
Assistant's Solution 2: <solution2>

Both solutions are correct. You MUST compare them based on the following criteria:

- The reasoning process is more correct, and logical.

- The reasoning process does not skip any reasoning steps.

- The reasoning process does not skip any reasoning steps.

You MUST follow this exact format:
Your detailed verification reasoning goes here. Conclude with the number of the preferred solution: $\boxed{1}$ or $\boxed{2}$.
If you prefer solution 1, you MUST output $\boxed{1}$.
If you prefer solution 2, you MUST output $\boxed{2}$.

Your preference:

---

## B ADDTIONAL EXPERIMENTAL RESULTS

In this section, we include additional ablation studies and evaluation results for a more comprehensive understanding of the PROF-GRPO framework. We conduct the following variants of PROF framework:

1. Ratio: remove correct and incorrect responses according to the original correct–incorrect ratio rather than enforcing a balanced ratio in PROF.

2. Filter-Nstep: Ranking and filtering out the samples with smaller number of steps instead of lower PRM-ORM consistency.

From Table 5, we find that Ratio scores 51.7% on average and cannot compete with balancing the proportion (PROF-GRPO), which also corroborates the conclusion that maintaining a balanced correct-incorrect proportion is essential. Additionally, since we observe that PROF boosts the number of intermediate reasoning steps, to verify that PROF does not simply increase the step length, but more importantly, enhances the quality of reasoning steps, we simply use the step length as the filtering criterion for comparison. As shown in Figure 7 and Table 5, Filter-Nstep manipulates the step length, which exhibits an unreasonable increase followed by a sudden drop, and its average accuracy is inferior.

## C ADDITIONAL EXAMPLES

| Algorithm | Math500 | Minerva Math | Olympiad Bench | AIME24 | AMC23 | Average |
|---|---|---|---|---|---|---|
| PROF-GRPO | **83.1** | **39.0** | **47.8** | 17.5 | **70.9** | **51.7** |
| Filter-Nstep | 81.5 | 35.5 | 45.9 | 16.3 | 58.6 | 47.6 |

Table 5: Performance of variants of filtering methods besides PROF-GRPO starting from Qwen2.5-Math-7B-base averaged on all five benchmarks, where Ratio denotes filtering while preserving the original correct-/incorrect proportion, Filter-Nstep denotes ranking and filtering based on the number of step segments.
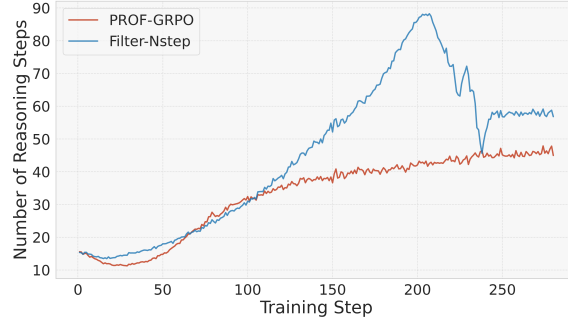


Figure 7: The number of reasoning steps during training time for PROF-GRPO and Filter-Nstep initialized from Qwen2.5-Math-7B-base.

**Problem:** The proper divisors of 12 are 1, 2, 3, 4 and 6. A proper divisor of an integer $N$ is a positive divisor of $N$ that is less than $N$. What is the sum of the proper divisors of the sum of the proper divisors of 284?

**PROF-GRPO Solution:**

First, we need to find the proper divisors of 284 …
Let's start by finding the divisors of 284. We can start by checking the obvious small divisors.
Let's check if 2 is a divisor of 284.
$$284 \div 2 = 142$$
Since the result is an integer, 2 is a divisor.
Next, let's check if 4 is a divisor of 284…
The proper divisors of 284 are 1, 2, 4, 71, and 142.
Next, we need to find the sum of these proper divisors. Let's add these step by step
…
$$1 + 2 + 4 + 71 + 142 = 220$$
Now, we need to find the proper divisors of 220. …
The proper divisors of 220 are 1, 2, 4, 5, 10, 11, 20, 22, 44, 55, and 110. The sum is:
$$1 + 2 + \ldots + 110 = 284$$
Therefore, the sum of the proper divisors of the sum of the proper divisors of 284 is $\boxed{284}$.

**GRPO Solution:**

First, we need to find the proper divisors of 284. The proper divisors of 284 are the positive divisors of 284 that are less than 284. The divisors of 284 are 1, 2, 4, 71, and 142. The sum of these proper divisors is:
$$1 + 2 + 4 + 71 + 142 = 220$$
Next, we need to find the proper divisors of 220. The divisors of 220 are 1, 2, 4, 5, 10, 11, 20, 22, 44, 55, and 110. The sum of these proper divisors is:
$$1 + 2 + \ldots + 110 = \boxed{284}$$
So the sum of the proper divisors of the sum of the proper divisors of 284 is 284 .

**Blend-PRM-GRPO Solution:**

First, we need to find the proper divisors of 284.
Let's start by finding all the divisors of 284.
We'll start with 1, which is always a divisor.
We check if 2 is a divisor.
$$284 \div 2 = 142$$
So, 2 is a divisor.
Next, we check if 3 is a divisor.
$$284 \div 3 \approx 94.67$$
So, 3 is not a divisor …
We continue this process until we reach 142.
The divisors of 284 are 1, 2, 4, 71, 142.
Now, we need to find the sum of these proper divisors: $1 + 2 + 4 + 71 + 142$
First, we add 1 and 2: $1 + 2 = 3$. ….
So, the sum of the proper divisors of 284 is 220.

Now, we need to find the proper divisors of 220.
We'll check each number from 1 to 219 to see if it divides 220 without leaving a remainder …
The proper divisors of 220 are 1, 2, 4, 5, 10, 11, 20, 22, 44, 55, and 110….
So, the sum of the proper divisors of 220 is $\boxed{284}$ .

Figure 8: A Math500 example to compare distinct intermediate reasoning patterns of PROF-GRPO, vanilla GRPO and Blend-PRM-GRPO. PROF-GRPO presents concrete and correct deduction steps. PROF-GRPO's solution shows how to find the divisors and summation in detail, and is easy to follow. GRPO skips all core reasoning. Blend-PRM-GRPO has inefficient and excessively tedious steps.