# GDPO: Group reward-Decoupled Normalization Policy Optimization for Multi-reward RL Optimization

**Shih-Yang Liu[1], Xin Dong[\*], Ximing Lu, Shizhe Diao, Peter Belcak, Mingjie Liu, Min-Hung Chen, Hongxu Yin, Yu-Chiang Frank Wang, Kwang-Ting Cheng[1], Yejin Choi, Jan Kautz, Pavlo Molchanov**

**NVIDIA**

As language models become increasingly capable, users expect them to provide not only accurate responses but also behaviors aligned with diverse human preferences across a variety of scenarios. To achieve this, Reinforcement learning (RL) pipelines have begun incorporating multiple rewards, each capturing a distinct preference, to guide models toward these desired behaviors. However, recent work has defaulted to apply Group Relative Policy Optimization (GRPO) under multi-reward setting without examining its suitability. In this paper, we demonstrate that directly applying GRPO to normalize distinct rollout reward combinations causes them to collapse into identical advantage values, reducing the resolution of the training signal and resulting in suboptimal convergence and, in some cases, early training failure. We then introduce **G**roup reward-**D**ecoupled Normalization **P**olicy **O**ptimization (**GDPO**), a new policy optimization method to resolve these issues by decoupling the normalization of individual rewards, more faithfully preserving their relative differences and enabling more accurate multi-reward optimization, along with substantially improved training stability. We compare GDPO with GRPO across three tasks: tool calling, math reasoning, and coding reasoning, evaluating both correctness metrics (accuracy, bug ratio) and constraint adherence metrics (format, length). Across all settings, GDPO consistently outperforms GRPO, demonstrating its effectiveness and generalizability for multi-reward reinforcement learning optimization.

**Implementations: HF-TRL, verl, Nemo-RL | Links: Project, Lab**



(a) An overview of our proposed GDPO
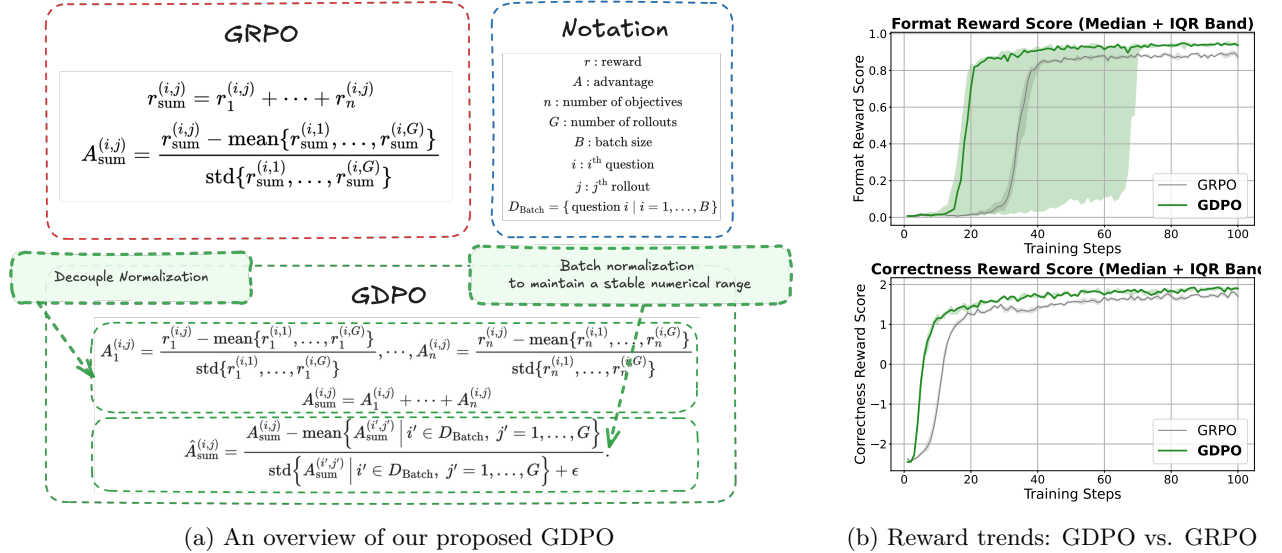
(b) Reward trends: GDPO vs. GRPO

Figure 1 | (a): An overview of GDPO, which performs group-wise normalization per reward and then applies batch-wise advantage normalization to preserve a stable numerical range independent of reward count and improve update stability. (b): Median and IQR reward curves over five runs of Qwen2.5-Instruct-1.5B tool-calling RL, demonstrating that GDPO consistently converges to higher correctness and format reward score than GRPO.

# 1. Introduction

As language models continue to advance in capability, expectations for their behavior have grown accordingly. Demand for models to not only provide accurate responses but also exhibit behaviors aligned with a wide range of human preferences across diverse scenarios has continued to increase. These preferences span efficiency [1, 2, 3], safety [4], response coherence and logic [5, 6], gender biases [7] and many other objectives. Meeting such heterogeneous requirements within a single model is a challenging task.

Reinforcement learning (RL) has emerged as the de facto training pipeline for aligning large language models to fulfill such diverse human preferences. In particular, recent RL-based approaches have begun to incorporate multiple rewards into training, with each reward designed to capture different human preferences and collectively guide models toward human-favored behaviors. Despite this growing interest in multi-reward RL, recent work [1, 3, 5] has largely focused on the reward design itself and often directly relied on applying Group Relative Policy Optimization (GRPO) directly for multi-reward RL optimization, often without examining whether GRPO is well-suited for optimizing combinations of heterogeneous rewards.

In this paper, we revisit the applicability of GRPO in multi-reward settings and show that directly applying GRPO to normalize different combinations of rollout rewards can cause them to collapse into identical advantage values, which effectively limits the precision of the training signal, as illustrated in Fig. 2. This collapse removes important distinctions across reward dimensions and leads to inaccurate policy updates, suboptimal reward convergence, and, in many cases, early training failure.

To overcome these challenges, we propose **Group reward-Decoupled Normalization Policy Optimization (GDPO)** which decouples the group-wise normalization of each individual reward as illustrated in Fig. 1a, to ensure that distinctions across different reward combinations are better preserved and more accurately reflect the relative differences in model responses. This leads to more precise multi-reward optimization and substantially improved training convergence. After this decoupled group-wise normalization, we apply batch-wise advantage normalization to ensure that the magnitude of advantage does not increase as the number of individual rewards increases.

We compare GDPO and GRPO across three tasks: tool calling, math reasoning, and code reasoning. These tasks cover a wide range of objectives, including tool-calling accuracy and format correctness, mathematical reasoning accuracy and adherence to reasoning-length constraints, and code pass rate and bug ratio. Across all tasks, GDPO converges better. For example, in Fig. 1b, training Qwen2.5-1.5B-Instruct with GDPO attains both higher correctness and format compliance than GRPO on the tool-calling task. On challenging math tasks, GDPO consistently outperforms GRPO. For instance, training DeepSeek-R1-1.5B and Qwen3-4B-Instruct with GDPO yields up to 6.3% and 2.3% higher accuracy on AIME compared to GRPO, while keeping more responses short simultaneously.

Taken together, these results demonstrate the effectiveness and generalizability of GDPO, showing it to be a better alternative to GRPO for multi-reward RL optimization.

Our contributions are as follows:

- **Analysis of GRPO reward collapse.** We demonstrate that applying GRPO naively for multi-reward RL optimization can collapse distinct rollout reward combinations into identical advantage values, thereby diminishing the resolution of the learning signal.

- **Remediation of GRPO reward collapse.** We propose GDPO, which performs group-wise decoupled normalization of each reward separately to better preserve cross-reward distinctions and enable more accurate multi-reward optimization.

- In addition to GDPO, we provide a systematic overview of how to modify reward functions and adjust reward weights to more faithfully align with preferences of varying priority.

- We carry out extensive experiments on three tasks: tool calling, math reasoning, and code reasoning, and compare the effectiveness of GDPO on optimizing a wide range of rewards corresponding to accuracy, format correctness, length constraints, and code quality. In all settings, GDPO consistently outperforms GRPO, showing improved training convergence and stronger downstream performance that align more closely with a diverse set of preferences.

## 2. GRPO's propensity for reward signal collapse in multi-reward RL

Recent advancements such as Group Relative Policy Optimization (GRPO) [8] and its variants, including DAPO [9] and Reinforce++-Baseline [10], have emerged as widely adopted reinforcement learning algorithms due to their efficiency and simplicity. In contrast to Proximal Policy Optimization (PPO) [11], GRPO eliminates the need for a value model by leveraging group-relative advantage estimation for policy updates.

Currently, GRPO has been primarily employed for optimizing a single-objective reward, typically focusing on accuracy. However, as model capability continues to grow, recent works have increasingly sought to optimize multiple rewards, such as response length constraint and formatting quality, in addition to accuracy [1, 12, 3], to better align with human preferences. Existing approaches for multi-reward RL generally adopt a straightforward strategy: summing all reward components and applying GRPO directly.

Formally, for a given question–answer pair $(q_i, o_j)$, where the behavior policy $\pi_{\theta_{\text{old}}}$ samples a group of $G$ responses $\{o_j\}_{j=1}^G$, and assuming $n$ objectives, the aggregated reward for the $j$-th response is computed as the sum of each objective's reward:

$$r_{\text{sum}}^{(i,j)} = r_1^{(i,j)} + \cdots + r_n^{(i,j)} \tag{1}$$

The group-relative advantage for the $j$-th response is then obtained by normalizing the group-level aggregated rewards:

$$A_{\text{sum}}^{(i,j)} = \frac{r_{\text{sum}}^{(i,j)} - \text{mean}\{r_{\text{sum}}^{(i,1)}, \ldots, r_{\text{sum}}^{(i,G)}\}}{\text{std}\{r_{\text{sum}}^{(i,1)}, \ldots, r_{\text{sum}}^{(i,G)}\}} \tag{2}$$

The corresponding multi-reward GRPO optimization objective can then be expressed as:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{(q_i, o_j) \sim D, \{o_j\}_{j=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|q)} \left[ \frac{1}{G} \sum_{j=1}^G \frac{1}{|o_j|} \sum_{t=1}^{|o_j|} \min\left( s_{i,t}(\theta) A_{\text{sum}}^{(i,j)}, \text{clip}(s_{i,t}(\theta), 1-\epsilon, 1+\epsilon) A_{\text{sum}}^{(i,j)} \right) \right] \tag{3}$$

where $s_t(\theta) = \frac{\pi_\theta(o_j^t \mid q, o_j^{<t})}{\pi_{\theta_{\text{old}}}(o_j^t \mid q, o_j^{<t})}$ and $\epsilon$ denotes the clipping threshold. For clarity, we omit the KL-divergence loss term in this formulation.

We first revisit this common practice of applying GRPO for mulit-reward RL optimization and identify a previously overlooked issue, that is GRPO inherently compresses the reward signal, causing loss of information in the advantage estimates. To illustrate, we start with a simple training setting and then extend it to more general cases. Consider a scenario where we generate two rollouts for each question for calculating the group-relative advantage and the task involves two binary reward $r_1, r_2 \in \{0, 1\}$. Consequently, the total reward for each rollout can take values from $\{0, 1, 2\}$.

As shown in Figure 2, we enumerate all possible rollout reward combinations within a group, represented as (`rollout 1's total reward`, `rollout 2's total reward`) and the corresponding normalized advantages as (*rollout 1's normalized advantage*, *rollout 2's normalized advantage*). Despite having six distinct combinations when order is ignored, only two unique advantage groups emerge after applying group-wise reward normalization. Specifically, (0, 1), (0, 2), and (1, 2) yield identical normalized advantages $A_{\text{sum}}$ of (−0.7071, 0.7071), while (0, 0), (1, 1), and (2, 2) all result in (0, 0).

This demonstrates a fundamental limitation of GRPO's advantage calculation in multi-reward optimization which over-compresses the rich group-wise reward signal. Intuitively, (0, 2) should produce a stronger learning signal than (0, 1) because a total reward of 2 indicates simultaneous satisfaction of two rewards, whereas a reward of 1 corresponds to achieving only one. Thus, when the other one rollout only receives zero reward, (0, 2) should yield a larger relative advantage than (0, 1). This limitation can also introduce risks of training instability due to inaccurate advantage estimates. As shown in Fig. 5, the correctness reward score begins to decline after approximately 400 training steps when training with GRPO, indicating a partial training collapse.

Recently, Dr.GRPO [13] and DeepSeek-v3.2 [6] adopt a variant of GRPO that removes the standard deviation normalization term from Eq. 2, such that $A_{\text{sum}}^{(i,j)} = r_{\text{sum}}^{(i,j)} - \text{mean}\{r_{\text{sum}}^{(i,1)}, \ldots, r_{\text{sum}}^{(i,G)}\}$. Despite these works introduce this modification to mitigate question-level difficulty bias, at first glance, this change also appears to
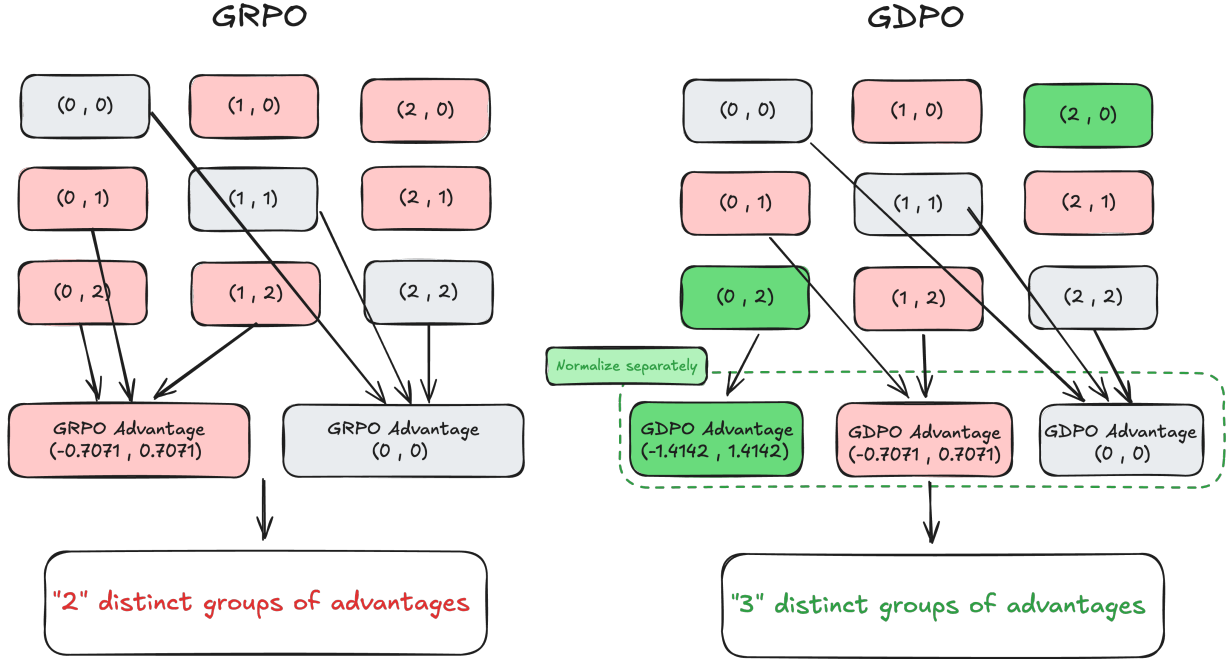
Figure 2 | Comparison of GRPO and GDPO advantage computation in a two-binary-reward, two-rollout example. GRPO maps different reward combinations into only two distinct advantage groups, whereas GDPO normalizes each reward independently and retains three distinct groups of advantage values. We skip the batch-wise normalization calculation step in GDPO here for simplicity since it does not change the number of distinct advantage groups.

address the issue we identify. Specifically, removing the standard deviation normalization mitigates the issue: $(0, 1)$ and $(0, 2)$ now yield distinct advantages of $(-0.5, 0.5)$ and $(-1.0, 1.0)$, respectively. However, when this setup is generalized to a larger number of rollouts while keeping the number of rewards fixed, as shown in Figure 3, we observe that such fix only slightly increases the number of distinct advantage groups compared to GRPO. A similar trend can be observed under settings where the number of rollouts is fixed at four, but the number of rewards gradually increases. In this case, we also observe only modest improvements in the number of distinct advantage groups. We also empirically examine the effectiveness of removing the standard deviation normalization term in Section 4.1.1, and find that this modification does not lead to improved convergence or better downstream evaluation performance.

## 3. Method

### 3.1. Group reward-Decoupled normalization Policy Optimization

To overcome these challenges, we propose **Group reward-Decoupled normalization Policy Optimization (GDPO)**, a method designed to better maintain distinctions among different reward combinations and more accurately capture their relative differences in the final advantages. In contrast to GRPO, which applies group-wise normalization directly to the aggregated reward sum, GDPO decouples this process by performing group-wise normalization of each reward separately before aggregation. Concretely, rather than summing all $n$ rewards first (as in Eq. 1) and then applying group-wise normalization to obtain $A_{\mathrm{sum}}$ (Eq. 2), GDPO computes the normalized advantage for each reward for the $j^{\mathrm{th}}$ rollout of the $i^{\mathrm{th}}$ question as:

$$A_1^{(i,j)} = \frac{r_1^{(i,j)} - \mathrm{mean}\{r_1^{(i,1)}, \ldots, r_1^{(i,G)}\}}{\mathrm{std}\{r_1^{(i,1)}, \ldots, r_1^{(i,G)}\}}, \quad \ldots, \quad A_n^{(i,j)} = \frac{r_n^{(i,j)} - \mathrm{mean}\{r_n^{(i,1)}, \ldots, r_n^{(i,G)}\}}{\mathrm{std}\{r_n^{(i,1)}, \ldots, r_n^{(i,G)}\}} \qquad (4)$$
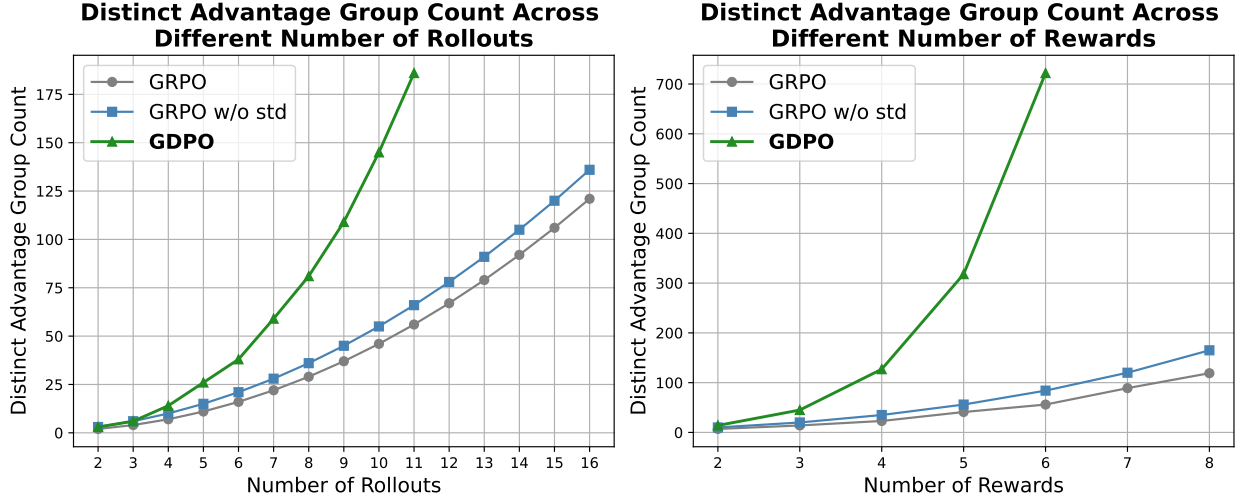
Figure 3 | Comparison of the number of distinct advantage groups produced by GRPO, GRPO without standard deviation normalization (GRPO w/o std), and GDPO. As the number of rollouts (left) or rewards (right) grows, GDPO consistently preserve a substantially larger number of distinct advantage groups compared to GRPO and GRPO w/o std. This results in advantage estimations that provide more expressive training signals.

The overall advantage used for policy updates is then obtained by first summing the normalized advantages across all objectives:

$$A_{\text{sum}}^{(i,j)} = A_1^{(i,j)} + \cdots + A_n^{(i,j)} \tag{5}$$

$$\hat{A}_{\text{sum}}^{(i,j)} = \frac{A_{\text{sum}}^{(i,j)} - \text{mean}\Big\{ A_{\text{sum}}^{(i',j')} \mid i' \in D_{\text{Batch}}, \ j' = 1, \ldots, G \Big\}}{\text{std}\Big\{ A_{\text{sum}}^{(i',j')} \mid i' \in D_{\text{Batch}}, \ j' = 1, \ldots, G \Big\} + \epsilon} \tag{6}$$

then applying batch-wise advantages normalization to the sum of the multi-reward advantages, which ensures that the numerical scale of the final advantage $\hat{A}_{\text{sum}}^{(i,j)}$ remains stable and does not grow as additional rewards are introduced. Empirically, we also find that this normalization step improves training stability, as shown in Appendix A, where removing batch-wise normalization occasionally leads to convergence failures.

By separating the normalization of each reward, GDPO alleviates the information-loss problem present in GRPO's advantage estimation, as illustrated in Fig. 2. Note that since the batch-wise normalization step in GDPO does not alter the number of distinct advantage groups, we omit it here for clarity. From the figure, we can see that when adopting GRPO, distinct reward combinations, such as $(0,2)$ and $(0,1)$, lead to identical normalized advantages, masking the subtle distinctions between them. In contrast, GDPO retains these fine-grained differences by assigning distinct advantage values to each combination, for example, the reward combination of $(0,1)$ after GDPO normalization becomes $(-0.7071, 0.7071)$ and $(0,2)$ becomes $(-1.4142, 1.4142)$, which more appropriately reflects that $(0,2)$ should yield a stronger learning signal than $(0,1)$. Similarly, when extending the number of rollouts to three, GRPO would assign advantage values of $(0,0,0)$ to $(1,1,1)$. However, $(1,1,1)$ may arise from heterogeneous reward partitions such as $r_1 = (1,1,0)$ or $r_2 = (0,0,1)$, for which GDPO would yield non-zero advantages, thereby preserving meaningful differences across reward dimension.

We further quantify the effectiveness of GDPO by comparing the number of distinct advantage groups across GDPO, GRPO, and GRPO w/o std under two experimental settings as shown in Fig. 3. In the two-reward scenario with a varying number of rollouts, GDPO consistently produces a significantly higher count of distinct advantage groups, with the gap widening as the number of rollouts increases. On the other hand, when fixing the number of rollouts to four and increasing the number of rewards, a similar pattern emerges, where GDPO exhibits progressively larger advantage granularity as the objective count grows. This demonstrate that the

proposed decoupled normalization approach effectively increases the number of distinct advantage groups across all the RL settings and enables more precise advantage estimation. In addition to these theoretical improvements, we observe that using GDPO consistently yields a more stable training curve and improved convergence. For instance, GDPO achieves better convergence on both the format reward and the correctness reward in the tool-calling task, as shown in Fig. 4. GDPO also eliminates the training collapse issue observed with GRPO in the math reasoning task, as shown in Fig. 5, where the model trained with GDPO continues to improve the correctness reward score throughout training. Additional empirical results in Section 4 further confirm GDPO's ability to achieve stronger alignment with the target preferences across a wide range of downstream tasks.

### 3.2. Effective incorporation of priority variation

Up to this point, we have assumed that all objectives carry equal importance. In practice, this assumption does not always hold in real-world applications. In this section, we provide a systematic overview of how to adjust the weights of rewards associated with different objectives or modify the reward functions to enforce prioritization of more important objectives. We also discuss how these two design choices behave differently when the underlying rewards vary significantly in difficulty.

It is common practice to assign different weights to each reward to encode different priorities among objectives, such that $r_{\text{sum}} = w_1 r_1 + \cdots + w_n r_n$, thereby controlling the contribution of each reward to the final advantage used for policy updates, and for GDPO, such weightings are apply to the normalized advantages of each reward as:

$$A_{\text{sum}}^{(i,j)} = w_1 A_1^{(i,j)} + \cdots + w_n A_n^{(i,j)} \tag{7}$$

However, we found that adjusting reward weights does not always yield the intended behavior when the difficulty levels of the underlying objectives differ substantially. If one objective is much easier than the others, the model often focuses on maximizing the reward for that objective regardless of the assigned weights. As a result, to more effectively force the model to allocate more attention to rewards associated with more challenging objectives, the weight differences must be made sufficiently large to compensate for the disparity in difficulty. Nevertheless, even with such adjustments, the model may still prefer optimizing the easier reward rather than the objective the user intends to prioritize, a phenomenon that we empirically demonstrate in Sec. 4.2.1.

Therefore, some recent works [1, 14] address such reward hacking by conditioning easier rewards on more difficult rewards. Specifically, Given two rewards $r_k$ and $r_l$, conditioning $r_k$ on $r_l$ can be formulated as:

$$r_k = \begin{cases} r_k, & \text{if } r_l \geq t \\ 0, & \text{otherwise.} \end{cases} \tag{8}$$

With such reward function design, the model can only receive reward of $r_k$ when the reward $r_l$ satisfies a predefined score threshold $t$, and as a consequence, the model get force to always maximize the human-prioritized reward first and completely alleviate the above mentioned problem. The empirical effectiveness of this strategy is shown in Sec. 4.2.1, where models trained with conditioned reward functions achieve higher performance on prioritized objectives compared to those trained without conditioning but only with larger weight assigned to the prioritized rewards. We also observe that, after resolving the issue of the easier reward dominating, assigning different reward weights for fine-grained priority adjustment can also be more faithfully reflected in the final model behavior.

## 4. Experiments

We begin by evaluating the effectiveness of GDPO compared with GRPO on the tool-calling task (Sec. 4.1), which involves optimizing two rewards: tool-calling correctness and format compliance. We then present an ablation study that examines the training convergence and downstream performance of GRPO with and without standard deviation normalization. Next, we compare GDPO and GRPO on a math reasoning task that optimizes two implicitly competing rewards, accuracy and length constraint (Sec. 4.2). We further conduct extensive analyses of the impact of incorporating different reward weights and modifying reward functions to better reflect varying priorities in human preferences, especially when rewards differ substantially

in difficulty. Finally, we extend the number of optimized rewards to three and compare GRPO and GDPO on coding reasoning (Sec. 4.3), jointly optimizing code-generation accuracy, adherence to length constraints, and bug ratio, further demonstrating that GDPO generalizes effectively to settings with three reward objectives.
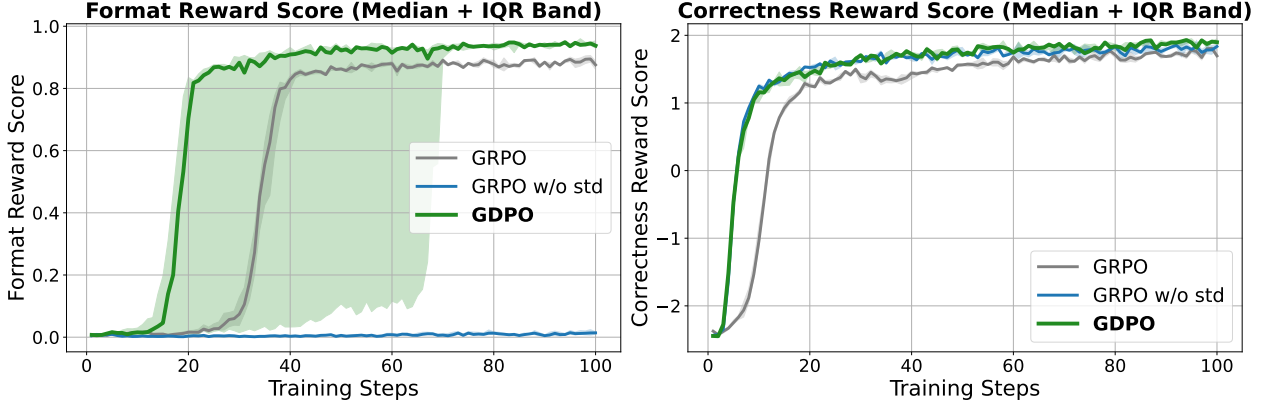
## 4.1. Tool calling



Figure 4 | Median and IQR reward curves across five runs of Qwen2.5-1.5B on the tool-calling task for GDPO, GRPO, and GRPO w/o std. GDPO consistently converges to higher correctness and format rewards, while GRPO w/o std matches correctness gains but fails to converge on the format reward.

We compare GDPO with GRPO on the tool calling task following the setup of ToolRL [12][1]. Specifically, the model is trained to learn how to incorporate external tools into the reasoning trajectory to solve a user task following the output format shown in Appendix B, where the reasoning steps must be enclosed in `<think></think>`, the tool calls must appear within `<tool_call></tool_call>`, and the model's final answer must be placed inside `<response></response>`. We adopt the same training set as ToolRL for RL training, which consists of 2k samples from ToolACE [15], 1k samples from Hammar [16], and 1k samples from xLAM [17]. Each training instance contains a question and its corresponding ground-truth tool calls. The training involves two rewards:

- **Format reward**: The format reward $\mathcal{R}_{\text{format}} \in \{0, 1\}$ checks whether the model output satisfies the required structure and contains all necessary fields in the correct order.

- **Correctness reward**: The correctness reward $\mathcal{R}_{\text{correct}} \in [-3, 3]$ evaluates the model-generated tool calls against the ground-truth calls using three metrics: tool name matching, parameter name matching, and parameter content matching.

A full description of the reward formulation is provided in Appendix C. We train Qwen-2.5-Instruct (1.5B and 3B) [18] with GRPO and GDPO using verl [19] for 100 steps, following the original hyperparameter settings from ToolRL's GRPO recipe. We use four rollouts per training question, a batch size of 512, and a maximum response length of 1024. The complete hyperparameter configuration is listed in Appendix D.

We evaluate the trained models on the Berkeley Function Call Leaderboard (BFCL-v3) [20], a comprehensive benchmark covering a broad range of challenges, including single-step reasoning, multi-step tool use, real-time execution, irrelevant tool rejection, simultaneous multi-tool selection, and multi-tool execution. We finetune the models with GRPO and GDPO across five runs and report the average accuracy and average format correctness on BFCL-v3 in Table 1. We additionally plot the median training curves with interquartile ranges for both methods over the five runs in Fig. 4.

From the training curves, we observe that GDPO consistently converges to higher values on both the format and correctness reward score across all runs. Although GDPO exhibits larger variance in the number of steps

---

[1]https://github.com/qiancheng0/ToolRL

required to converge on format reward, it ultimately attains better format compliance than both GRPO. For the correctness reward, GDPO shows faster early-stage improvement and reaches a higher reward score than the GRPO baselines toward later stages, demonstrating the effectiveness of GDPO on providing more accurate advantage estimation that leads to better optimization.

Table 1 | Comparison of GDPO and GRPO-trained Qwen2.5-Instruct-1.5B/3B models on tool-calling accuracy and format correctness. The reported results are the averages across five runs.

| | Live Ocerall Acc ↑ | Multi Turn Overall Acc ↑ | Non-Live Overall Acc ↑ | Avg Acc ↑ | Correct Format ↑ |
|---|---|---|---|---|---|
| Qwen2.5-Instruct-1.5B | 37.89% | 0.12% | 15.63% | 17.88% | 4.74% |
| GRPO | 50.63% | 2.04% | 37.87% | 30.18% | 76.33% |
| GDPO | **55.36%** | **2.50%** | **40.58%** | **32.81%** | **80.66%** |
| Qwen2.5-Instruct-3B | 63.57% | 1.38% | 30.75% | 31.90% | 58.37% |
| GRPO | 69.23% | 3.14% | 45.24% | 39.20% | 81.64% |
| GDPO | **71.22%** | **4.59%** | **46.79%** | **40.87%** | **82.23%** |

In the BFCL-v3 evaluation shown in Table 1, GDPO also consistently improves the average tool calling accuracy and format correctness over the GRPO-trained counterparts. For training Qwen2.5-Instruct-1.5B, GDPO achieves almost 5% and 3% improvement on Live/non-Live tasks and gains roughly 2.7% improvement on the overall average accuracy and more than 4% in correct format ratio compared with GRPO. Similar improvements are observed for the 3B model, where GDPO continues to outperform GRPO across all the sub-tasks, achieving up to 2% accuracy improvement and delivers a better format compliance ratio.

### 4.1.1. Does removing the standard deviation normalization term in GRPO provide any benefit?

Table 2 | Comparison of GRPO, GRPO w/o std, and GDPO-trained Qwen2.5-Instruct-1.5B/3B models on tool-calling accuracy and format correctness.The reported results are the average across five runs.

| | Live Ocerall Acc ↑ | Multi Turn Overall Acc ↑ | Non-Live Overall Acc ↑ | Avg Acc ↑ | Correct Format ↑ |
|---|---|---|---|---|---|
| Qwen2.5-1.5B-Instruct | 37.89% | 0.12% | 15.63% | 17.88% | 4.74% |
| GRPO | 50.63% | 2.04% | 37.87% | 30.18% | 76.33% |
| GRPO w/o std | 47.19% | 1.47% | 39.11% | 29.26% | 0% |
| GDPO | **55.36%** | **2.50%** | **40.58%** | **32.81%** | **80.66%** |

Recall from Fig. 3 that removing the standard deviation normalization term in GRPO (denoted GRPO w/o std) slightly increases the number of distinct advantage groups. In this section, we empirically examine the effectiveness of this modification. Following the previous experiments, we run GRPO w/o std five times and report the average accuracy and average format correctness ratio on BFCL-v3.

In the reward training curves shown in Fig. 1b, we observe that although GRPO w/o std converges to a correctness reward that is similar to GDPO and higher than standard GRPO, it fails to improve the format reward entirely. This failure results in a correct format ratio of 0% on BFCL-v3 (see Table. 2), indicating that the model does not learn the required output structure. These also show that simply removing the standard deviation normalization term in order to increase advantage diversity can introduce instability into training, which may ultimately prevent successful convergence in multi-reward reinforcement learning.

### 4.2. Mathematical reasoning

We consider a mathematical reasoning task that optimizes two implicitly competing rewards: accuracy and adherence to a length constraint. The goal is to improve model performance on challenging mathematical problems while keeping the generated output within a predefined response length to encourage efficient problem solving. We train DeepSeek-R1-1.5B, DeepSeek-R1-7B [8], and Qwen3-4B-Instruct [21] using GRPO and GDPO on the DeepScaleR-Preview dataset [22] for 500 steps, which contains 40k competition-level math problems. Training is performed using verl [19], and we follow the original DeepSeek-R1 prompt format [8]. Following the DLER setup [14], we incorporate dynamic sampling, higher clipping thresholds, and the token-mean loss from DAPO [9], and use 16 rollouts, a batch size of 512, and a maximum response length of 8000 tokens. The full set of hyperparameters is provided in Appendix E.

The training uses two rewards:

- **Length reward:** The length reward $\mathcal{R}_{\text{length}} \in \{0, 1\}$ checks whether the model's output remains within the target length $l$, which is set to 4000 tokens for all remaining experiments:

$$\mathcal{R}_{\text{length}} = \begin{cases} 1, & \text{if response length} \leq l \\ 0, & \text{otherwise.} \end{cases}$$

- **Correctness reward:** The correctness reward $\mathcal{R}_{\text{correct}} \in \{0, 1\}$ indicates whether the final answer extracted from the model's response matches the ground truth.

We compare the GRPO and GDPO-trained model on AIME-24 [23], AMC (AMC 2022 and AMC 2023) [24], MATH [25], Minerva [26] and Olympiad Bench [27]. All evaluations are conducted using vLLM as the inference backend with a sampling temperature of 0.6, $top_p = 0.95$, and a maximum response length of 32k tokens. For each evaluation question, we generate 16 samples and report the average pass@1 score and the average length-exceeding ratio, denoted Exceed, which measures the percentage of model responses that exceed the predefined length limit of 4000 tokens.



Figure 5 | Training behavior of GRPO and GDPO on DeepSeek-R1-1.5B across correctness reward, length reward, and maximum batch response length. Both methods rapidly maximize the length reward, briefly suppressing correctness, yet GDPO subsequently recovers it and surpasses GRPO. After roughly 400 steps, GRPO's correctness score declines and its length-constraint violations increase, as reflected by rising maximum response lengths. In contrast, GDPO continues to improve correctness while steadily improving the control over response length.

From the training curves of GRPO and GDPO on DeepSeek-R1-1.5B as shown in Fig. 5, we first observe that the model tends to maximize the easier reward regardless of the optimization method. In this case, the length reward is easier to optimize, and both GRPO and GDPO reach a full length score within roughly the first 100 training steps. We also see that this rapid rise in the length reward coincides with an early drop in the correctness reward, which indicates that the two rewards are competing. During the initial phase of training, the model prioritizes satisfying the length constraint, often at the expense of the more challenging correctness objective. However, from the correctness reward trajectories, we observe that GDPO recovers the correctness reward more effectively than GRPO, achieving higher correctness scores at comparable training steps. We also see that GRPO training starts to destabilize after 400 steps with the correctness rewards score gradually decreasing while GDPO continue to improve the correctness score. Moreover, although both GDPO and GRPO maintain nearly perfect length scores throughout training, we also record the maximum response length within each training batch to assess how well the models satisfy the length constraint under more extreme cases. The results show that, despite achieving almost a full length reward, the maximum response length for GRPO begins to increase sharply after roughly 400 training steps, while the maximum response length for GDPO continues to decrease. Similar observation can be seen on the training curves on DeepSeek-R1-7B and Qwen3-4B-Instruct as shown in Fig 9 and Fig 10 in appendix where we can see that GDPO consistently provide better alignment to the length constraint. This contrast further illustrates the effectiveness of GDPO in multi-reward optimization compared with GRPO.

In addition, the benchmark results in Table 3 show that the GDPO-trained models not only achieve substantial improvements in reasoning efficiency over the original models, with up to a 80% reduction in length-exceeding

Table 3 | Comparison of GRPO and GDPO-trained DeepSeek-R1-1.5B/7B models on Pass@1 accuracy and the proportion of responses exceeding the length constraint across mathematical reasoning benchmarks.

| | | DeepSeek-R1-1.5B | | | DeepSeek-R1-7B | | | Qwen3-4B-Instruct | | |
| | | - | GRPO | GDPO | - | GRPO | GDPO | - | GRPO | GDPO |
|---|---|---|---|---|---|---|---|---|---|---|
| MATH | Acc ↑ | 84.3% | 83.6% | **86.2%** | 93.6% | **94.1%** | 93.9% | 94.6% | 93.9% | 93.9% |
| | Exceed ↓ | 35.0% | 1.5% | **0.8%** | 26.0% | 0.5% | **0.1%** | 11.3% | 0.8% | **0.1%** |
| AIME | Acc ↑ | 29.8% | 23.1% | **29.4%** | 55.4% | 50.2% | **53.1%** | 63.7% | 54.6% | **56.9%** |
| | Exceed ↓ | 91.5% | 10.8% | **6.5%** | 85.6% | 2.1% | **0.2%** | 71.3% | **2.5%** | 0.1% |
| AMC | Acc ↑ | 62.0% | 64.5% | **69.0%** | 82.9% | 83.8% | **84.0%** | 84.5% | **85.2%** | 84.3% |
| | Exceed ↓ | 67.5% | 3.2% | **2.3%** | 57.2% | 0.6% | **0.3%** | 33.9% | 0.7% | **0.1%** |
| Minerva | Acc ↑ | 38.41.% | 43.5% | **44.0%** | 49.8% | 53.2% | **53.8%** | 50.7% | **52.4%** | 51.9% |
| | Exceed ↓ | 51.4% | 1.7% | **0.3%** | 41.8% | 0.2% | **0.1%** | 9.1% | 0.3% | **0.1%** |
| Olympiad | Acc ↑ | 44.1% | 44.3% | **46.6%** | 58.2% | **60.2%** | 59.7% | 65.7% | 66.8% | **67.5%** |
| | Exceed ↓ | 70.1% | 2.6% | **1.9%** | 60.6% | 1.1% | **0.4%** | 41.3% | 1.6% | **1.0%** |

ratios on AIME, but also deliver higher accuracy on the majority of the tasks. Moreover, GDPO generally outperforms GRPO on both the accuracy and length constraint objectives. For the DeepSeek-R1-1.5B, GDPO outperforms GRPO across all benchmarks, achieving accuracy improvements of 2.6%/6.7%/2.3% on MATH, AIME and Olympiad, respectively, while also reducing the length exceed ratios across all the tasks. A similar trend holds for DeepSeek-R1-7B and Qwen3-4B-Instruct, where GDPO achieves stronger accuracy–efficiency trade-offs. The gains are particularly notable on the more challenging AIME benchmark, with GDPO improving accuracy by nearly 3% while reducing the length-exceeding rate to 0.2% and 0.1%, compared with 2.1% and 2.5% under GRPO for DeepSeek-R1-7B and Qwen3-4B-Instruct. Together, these results show that GDPO not only improves reasoning accuracy across a range of mathematical tasks but also adheres to the length constraint more effectively, underscoring its advantage in multi-reward optimization.

### 4.2.1. Impact analysis of different reward priority variation configurations
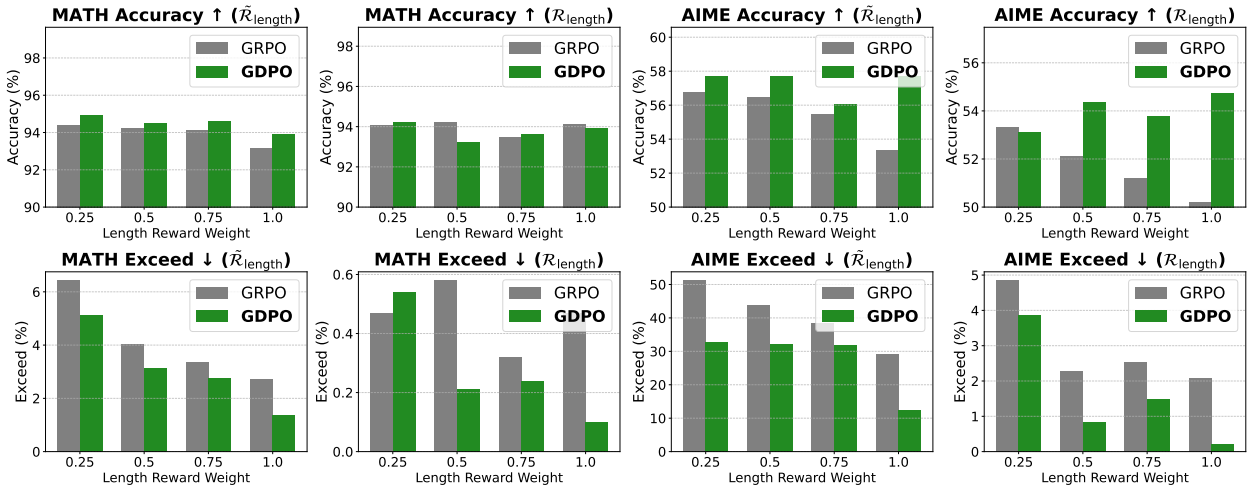


Figure 6 | Average accuracy and exceed-length ratios for GRPO/GDPO-trained DeepSeek-R1-7B models under varying length reward weights $\{1.0, 0.75, 0.5, 0.25\}$, with and without the conditioned length reward $\tilde{\mathcal{R}}_{\text{length}}$, on mathematical reasoning tasks.

Until this point, we have assumed that all rewards are treated with equal priority. However, as shown in

Fig. 5, the model often maximizes the easier objective at the cost of the more challenging one, even when both objectives are assigned the same reward weight. In this section, we investigate whether adjusting reward weights can guide the model to prioritize maximizing the correctness reward over the length reward when such a preference is desired and when the two objectives differ noticeably in difficulty.

We begin by fixing the reward weight for $\mathcal{R}_{\text{correct}}$, denoted $w_{\text{correct}}$, to 1, and varying the reward weight for $\mathcal{R}_{\text{length}}$, denoted $w_{\text{length}}$, over the set $\{0.25, 0.5, 0.75, 1.0\}$. This setup allows us to study whether reducing $w_{\text{length}}$ encourages the model to prioritize maximizing the more challenging correctness reward first. We carry out this experiment on DeepSeek-R1-7B and plot the average accuracy and average length-exceeding ratio of MATH and AIME in Fig. 6. Full results for the remaining tasks are provided in Appendix G.

From the results, we observe that reducing $w_{\text{length}}$ to 0.75 or 0.5 has little impact on the average length-exceeding ratio, which shifts by only 0.4% and 0.2% for GRPO on AIME and by 1.3% and 0.6% for GDPO. In addition, lowering $w_{\text{length}}$ does not necessarily relax the length constraint, as decreasing $w_{\text{length}}$ from 0.75 to 0.5 does not consistently increase the length-exceeding ratio on either AIME or MATH for GRPO or GDPO. This suggests that simply adjusting reward weights does not reliably induce the intended prioritization when the underlying objectives differ substantially in difficulty. Only when $w_{\text{length}}$ is reduced to 0.25, making it sufficiently small to compensate for the difficulty gap between the objectives, do we observe a clear increase in the length-exceeding ratio on AIME for both GRPO and GDPO and on MATH for GDPO.
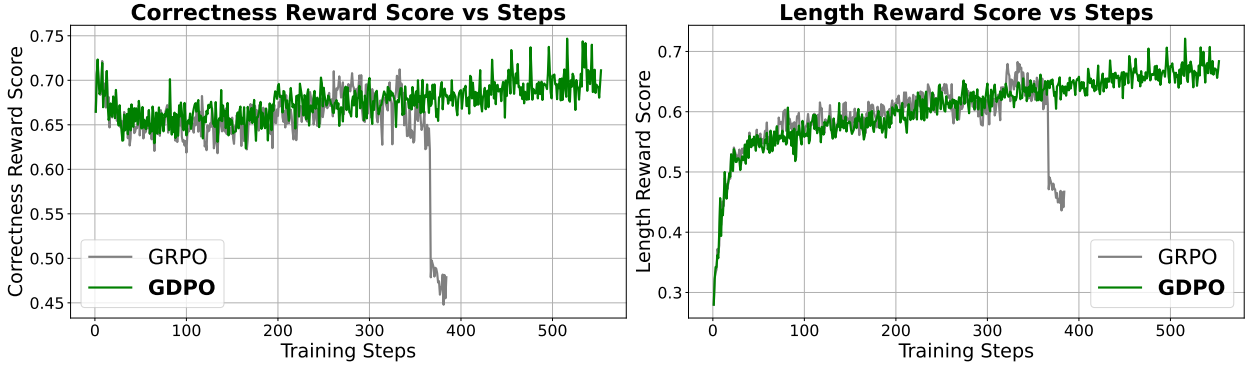


Figure 7 | Training curves of GRPO and GDPO with the conditioned length reward $\tilde{\mathcal{R}}_{\text{length}}$ on DeepSeek-R1-7B across correctness reward, length reward.

We next investigate whether conditioning the easier length reward on the more challenging correctness reward can help mitigate the disparity in difficulty between the two objectives and help improve priority alignment. Following the formulation in Sec. 3.2, we replace the original length reward $\mathcal{R}_{\text{length}}$ with a conditioned length reward defined as:

$$\tilde{\mathcal{R}}_{\text{length}} = \begin{cases} 1, & \text{if response length} \leq l \text{ and } \mathcal{R}_{\text{correct}} = 1 \\ 0, & \text{otherwise.} \end{cases}$$

Under this formulation, the model receives the length reward only when the generated response is also correct.

First, we observe that adopting the modified reward function $\tilde{\mathcal{R}}_{\text{length}}$ prevents the model from aggressively maximizing the length reward at the start of training. This reward design also helps avoid large drops in correctness reward score as the model tries to satisfy the length constraint. We can see that the average correctness reward decreases only slightly early in training and then gradually recovers from Fig. 7.

From Table 4, we also observe that using $\tilde{\mathcal{R}}_{\text{length}}$ leads to a larger increase in the average length-exceeding ratio for both GRPO and GDPO compared with merely adjusting the weight $w_{\text{length}}$ of $\mathcal{R}_{\text{length}}$, indicating a more effective relaxation of the length constraint. However, GRPO fails to convert this relaxed constraint into meaningful accuracy improvements. In contrast, GDPO prioritizes the correctness reward more effectively and achieves more consistent accuracy improvement over training without $\tilde{\mathcal{R}}_{\text{length}}$, while introducing substantially smaller increases in length violations. For instance, using $\tilde{\mathcal{R}}$length with GDPO yields a 4.4% accuracy improvement on AIME with a 16.9% reduction in length-exceeding ratio, and a 3% accuracy gain on AMC with a 4.8% reduction in length violations compared with using GRPO with the same reward.

Table 4 | Comparison of GRPO and GDPO DeepSeek-R1-7B models, with and without the conditioned length reward $\tilde{\mathcal{R}}_{\text{length}}$, on Pass@1 accuracy and the ratio of outputs exceeding the length constraint across mathematical reasoning benchmarks.

| | | DeepSeek-R1-7B | | | | |
|---|---|---|---|---|---|---|
| | | - | $\mathcal{R}_{\text{length}}$ | | $\tilde{\mathcal{R}}_{\text{length}}$ | |
| | | - | GRPO | GDPO | GRPO | GDPO |
| MATH | Acc ↑ | 93.6% | **94.1%** | 93.9% | 93.2% | **93.9%** |
| | Exceed ↓ | 26.0% | 0.5% | **0.1%** | 2.7% | **1.4%** |
| AIME | Acc ↑ | 55.4% | 50.2% | **53.1%** | 53.3% | **57.7%** |
| | Exceed ↓ | 85.6% | 2.1% | **0.2%** | 29.2% | **12.3%** |
| AMC | Acc ↑ | 82.9% | 83.8% | **84.0%** | 82.9% | **85.9%** |
| | Exceed ↓ | 57.2% | 0.6% | **0.3%** | 8.6% | **3.8%** |
| Minerva | Acc ↑ | 49.8% | 53.2% | **53.8%** | 53.2% | **53.4%** |
| | Exceed ↓ | 41.8% | 0.2% | **0.1%** | 2.5% | **1.0%** |
| Olympiad | Acc ↑ | 58.2% | **60.2%** | 59.7% | 59.1% | **60.8%** |
| | Exceed ↓ | 60.6% | 1.1% | **0.4%** | 10.3% | **6.2%** |

We next examine whether, after mitigating the difficulty disparity through conditioned length reward, varying the reward weight on $\tilde{\mathcal{R}}_{\text{length}}$, denoted $\tilde{w}_{\text{length}}$, leads to more faithful reflection of fine-grained preference adjustments. We fix the correctness reward weight and vary $\tilde{w}_{\text{length}} \in \{0.25, 0.5, 0.75, 1.0\}$. As shown in Fig. 6, the models trained with conditioned reward behave more predictably. For example, reducing $\tilde{w}_{\text{length}}$ from 1.0 to 0.25 steadily increases the length-exceeding ratio for both GRPO and GDPO on MATH and AIME, in contrast to the unstable results observed when adjusting the weight of the original $\mathcal{R}_{\text{length}}$.

Finally, across all settings, including different reward formulations and different reward weights, GDPO consistently provides a better accuracy and efficiency trade-off than GRPO.

### 4.3. Coding reasoning

Table 5 | Comparison of GRPO and GDPO trained DeepSeek-R1-7B models on coding pass rate, length-exceeding rate, and bug ratio across coding reasoning benchmarks. Here, $\mathcal{R}_{\text{pass}} + \tilde{\mathcal{R}}_{\text{length}}$ refers to optimizing $\mathcal{R}_{\text{pass}}$ and $\tilde{\mathcal{R}}_{\text{length}}$, and $\mathcal{R}_{\text{pass}} + \tilde{\mathcal{R}}_{\text{length}} + \mathcal{R}_{\text{bug}}$ refers to optimizing $\mathcal{R}_{\text{pass}}$, $\tilde{\mathcal{R}}_{\text{length}}$, and $\mathcal{R}_{\text{bug}}$.

| | | DeepSeek-R1-7B | | | | |
|---|---|---|---|---|---|---|
| | | - | $\mathcal{R}_{\text{Pass}} + \tilde{\mathcal{R}}_{\text{length}}$ | | $\mathcal{R}_{\text{Pass}} + \tilde{\mathcal{R}}_{\text{length}} + \mathcal{R}_{\text{Bug}}$ | |
| | | - | GRPO$_{\text{2-obj}}$ | GDPO$_{\text{2-obj}}$ | GRPO$_{\text{3-obj}}$ | GDPO$_{\text{3-obj}}$ |
| Apps | Pass ↑ | 28.1% | 67.2% | **68.3%** | **68.1%** | 67.8% |
| | Exceed ↓ | 73.9% | 5.2% | **5.0%** | 11.2% | **8.5%** |
| | Bug ↓ | 32.9% | 25.0% | **23.5%** | 20.3% | **18.8%** |
| Codecontests | Pass ↑ | 47.3% | 63.2% | **65.8%** | **65.6%** | **65.6%** |
| | Exceed ↓ | 83.0% | **14.2%** | 14.3% | 19.3% | **15.8%** |
| | Bug ↓ | 29.7% | 14.1% | **13.2%** | 3.9% | **2.5%** |
| Codeforces | Pass ↑ | 46.5% | 68.1% | **71.2%** | **69.5%** | 69.4% |
| | Exceed ↓ | 82.8% | **18.1%** | 18.4% | 16.9% | **13.6%** |
| | Bug ↓ | 27.8% | 7.0% | **5.6%** | 2.5% | **1.8%** |
| Taco | Pass ↑ | 28.1% | 45.1% | **48.4%** | 44.4% | **45.1%** |
| | Exceed ↓ | 78.0% | 11.8% | **10.8%** | 14.7% | **10.6%** |
| | Bug ↓ | 48.9% | 37.7% | **36.2%** | 30.0% | **28.0%** |

We examine whether GDPO continues to outperform GRPO when optimizing more than two rewards on our coding reasoning task. Similar to the mathematical reasoning setup, the objective is to improve the model's coding performance while constraining its output to a predefined target length. In addition, we introduce a third objective that encourages the model to generate bug-free code. We compare GDPO and GRPO by training DeepSeek-R1-7B on the Eurus-2-RL dataset [28], which contains 24k coding problems, each with multiple test cases. Training is conducted using the verl [19] framework for 400 steps, and we adopt the same hyperparameter configuration used in the mathematical reasoning experiments. The training optimizes three rewards:

- **Passrate reward:** The passrate reward $\mathcal{R}_{\text{pass}} \in [0, 1]$ measures the proportion of test cases passed by the generated code:
$$\mathcal{R}_{\text{pass}} = \frac{\text{number of passed test cases}}{\text{total test cases}}.$$

- **Conditioned Length reward:** The length reward $\tilde{\mathcal{R}}_{\text{length}} \in \{0, 1\}$ checks whether the model's response remains within the target length $l$ and whether the generated code satisfies correctness requirements:
$$\tilde{\mathcal{R}}_{\text{length}} = \begin{cases} 1, & \text{if response length } \leq l \text{ and } \mathcal{R}_{\text{pass}} = 1, \\ 0, & \text{otherwise.} \end{cases}$$

- **Bug reward:** The bug reward $\mathcal{R}_{\text{bug}} \in \{0, 1\}$ indicates whether the generated code runs without runtime or compilation errors.

For evaluation, we assess the trained model on the validation set from PRIME [28], which includes Apps [29], CodeContests [30], Codeforces[2], and Taco [31]. Following the same settings used for the mathematical reasoning evaluations, we use a sampling temperature of 0.6, a $top_p$ value of 0.95, and a maximum response length of 32k tokens. For each evaluation question, we generate 16 rollouts and report the average test case pass rate, the average length-exceeding ratio, and the average bug ratio, where the bug ratio measures the proportion of generated code that results in either a runtime error or a compilation error.

We compare GDPO and GRPO under two configurations: (1) a two-reward setting using $\mathcal{R}_{\text{pass}}$ and $\tilde{\mathcal{R}}_{\text{length}}$, and (2) a three-reward setting using $\mathcal{R}_{\text{pass}}$, $\tilde{\mathcal{R}}_{\text{length}}$, and $\mathcal{R}_{\text{bug}}$. We denote the two-reward and three-reward versions of GRPO as GRPO$_{\text{2-obj}}$ and GRPO$_{\text{3-obj}}$, and use the same notation for GDPO. As shown in Table 5, GDPO$_{\text{2-obj}}$ improves pass rates across all the tasks compared with GRPO$_{\text{2-obj}}$, while maintaining a similar length-exceeding ratio. For example, GDPO$_{\text{2-obj}}$ improves the Codecontests pass rate by 2.6% while increasing the length-exceeding ratio by only 0.1%, and achieves a 3.3% pass rate gain together with a 1% reduction in length violations compared with GRPO$_{\text{2-obj}}$ on Taco. A similar pattern holds in the three-reward setting where GDPO$_{\text{3-obj}}$ achieves a substantially better balance across all objectives, maintaining similar pass rate to GRPO$_{\text{3-obj}}$ while also markedly reducing both the length-exceeding ratio and the bug ratio.

Overall, these results demonstrate that GDPO remains effective as the number of reward signals increases. It consistently achieves a more favorable trade-off across all objectives than GRPO in both the two-reward and three-reward configurations.

## 5. Related Work

**GRPO Variants** Several extensions of Group Relative Policy Optimization (GRPO) [32] have been proposed to enhance the stability, effectiveness, and efficiency of the framework. These methods explore alternative formulations of group-wise normalization or policy updates while remaining grounded in the core GRPO principles. For example, to improve stability, Group Sequence Policy Optimization (GSPO) [33] defines the importance ratio based on sequence likelihood rather than at the token level, performing sequence-level clipping, rewarding, and optimization. To improve RL performance, Decoupled Clip and Dynamic sAmpling Policy Optimization (DAPO) [34] introduces four key techniques: Clip-Higher, Dynamic Sampling, Token-Level Policy Gradient Loss, and Overlong Reward Shaping. To promote efficient reasoning, Group Filtered Policy Optimization (GFPO) [35] addresses length explosion by sampling larger groups per

---

[2]https://huggingface.co/datasets/MatrixStudio/Codeforces-Python-Submissions

problem during training and filtering responses based on their length and reward-per-token ratio. Along the same direction, Doing Length pEnalty Right (DLER) [36] proposes a training recipe combining batch-wise reward normalization, higher clipping, dynamic sampling, and a simple truncation length penalty, achieving state-of-the-art accuracy–efficiency trade-offs.

**Multi-Reward Reinforcement Learning**  A growing body of work investigates RL approaches that incorporate multiple reward signals. One major usage is to model diverse human preferences. For example, Safe Reinforcement Learning from Human Feedback [37] decouples human preferences regarding helpfulness and harmlessness, dynamically adjusting the balance between the two objectives during fine-tuning. Similarly, Reinforcement Learning from Personalized Human Feedback (RLPHF) [38] optimizes LLMs for multiple (sometimes conflicting) preferences by training distinct policy models for each preference and merging them during inference. ALARM (Align Language Models via Hierarchical Rewards) [39] introduces a hierarchical reward structure that jointly captures dimensions such as response quality, style, fairness, and coherence. Recent developments in LLMs also integrate multiple-reward optimization to handle complex tasks. For instance, DeepSeek V3.2 [40] integrates rule-based outcome rewards, length penalties, and language-consistency rewards to enhance reasoning and agentic capabilities. Another important recent application for multi-reward RL is improving the efficiency of reasoning models while maintaining task performance, primarily by introducing length-based reward functions alongside outcome-based rewards. For example, O1-Pruner [41] and [42] apply normalized length penalties to ensure proportional compression. Similarly, [43] promotes conciseness by penalizing deviations from the shortest correct response within a sampled group. L1 [44] introduces Length Controlled Policy Optimization (LCPO) to optimize for accuracy while ensuring responses do not exceed a target length. Finally, [45] proposes an adaptive reward-shaping method that dynamically adjusts the trade-off between accuracy and response length based on model performance.

## 6. Conclusion

In contrast to prior work that focuses on designing new reward functions for multi-reward reinforcement learning while assuming GRPO is the default optimization method, this study revisits a fundamental but often overlooked question: whether GRPO is actually suitable for multi-reward optimization. Our analysis shows that applying GRPO directly to the summed reward can cause different reward combinations to collapse into the same advantage values. This collapse eliminates important distinctions across reward dimensions, produces inaccurate policy updates and weaker optimization performance, and can in many cases lead to early training failure.

To address this limitation, we introduce Group-wise Decoupled Policy Optimization (GDPO), a simple and effective modification to GRPO tailored for multi-reward reinforcement learning. GDPO performs normalization separately for each reward to preserve cross-reward differences, and it incorporates batch-wise advantage normalization to maintain a stable numerical range as additional rewards are included. These changes result in better convergence behavior and models that more faithfully reflect the intended preference structure.

We further present a systematic study for incorporating human preference priorities into the training process and explain how reward functions can be adjusted when the difficulty disparity between objectives is large. Through extensive experiments on tool calling, math reasoning, and coding reasoning, we show that GDPO consistently outperforms GRPO. Its advantages hold across different numbers of rewards, across different models, and across different reward functions.

Overall, our findings establish GDPO as a more stable, accurate, and preference-aligned optimization method than GRPO for multi-reward reinforcement learning, making it a strong foundation for aligning language models with diverse human preferences in real-world settings.

# References

[1] Wei Liu, Ruochen Zhou, Yiyun Deng, Yuzhen Huang, Junteng Liu, Yuntian Deng, Yizhe Zhang, and Junxian He. Learn to reason efficiently with adaptive length-based reward shaping. *arXiv preprint arXiv:2505.15612*, 2025.

[2] Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025.

[3] Pranjal Aggarwal and Sean Welleck. L1: Controlling how long a reasoning model thinks with reinforcement learning. *arXiv preprint arXiv:2503.04697*, 2025.

[4] Tong Mu, Alec Helyar, Johannes Heidecke, Joshua Achiam, Andrea Vallone, Ian Kivlichan, Molly Lin, Alex Beutel, John Schulman, and Lilian Weng. Rule based rewards for language model safety. *Advances in Neural Information Processing Systems*, 37:108877–108901, 2024.

[5] Yi Chen, Yuying Ge, Rui Wang, Yixiao Ge, Junhao Cheng, Ying Shan, and Xihui Liu. Grpo-care: Consistency-aware reinforcement learning for multimodal reasoning, 2025.

[6] Aixin Liu, Aoxue Mei, Bangcai Lin, Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao Wu, Bowei Zhang, Chaofan Lin, Chen Dong, et al. Deepseek-v3. 2: Pushing the frontier of open large language models. *arXiv preprint arXiv:2512.02556*, 2025.

[7] Tao Zhang, Ziqian Zeng, YuxiangXiao YuxiangXiao, Huiping Zhuang, Cen Chen, James R Foulds, and Shimei Pan. Genderalign: An alignment dataset for mitigating gender bias in large language models. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11293–11311, 2025.

[8] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

[9] 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*, 2025.

[10] Jian Hu, Jason Klein Liu, Haotian Xu, and Wei Shen. Reinforce++: An efficient rlhf algorithm with robustness to both prompt and reward models. *arXiv preprint arXiv:2501.03262*, 2025.

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

[12] Cheng Qian, Emre Can Acikgoz, Qi He, Hongru Wang, Xiusi Chen, Dilek Hakkani-Tür, Gokhan Tur, and Heng Ji. Toolrl: Reward is all tool learning needs. *arXiv preprint arXiv:2504.13958*, 2025.

[13] Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*, 2025.

[14] Shih-Yang Liu, Xin Dong, Ximing Lu, Shizhe Diao, Mingjie Liu, Min-Hung Chen, Hongxu Yin, Yu-Chiang Frank Wang, Kwang-Ting Cheng, Yejin Choi, et al. Dler: Doing length penalty right-incentivizing more intelligence per token via reinforcement learning. *arXiv preprint arXiv:2510.15110*, 2025.

[15] Weiwen Liu, Xu Huang, Xingshan Zeng, Xinlong Hao, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Zhengying Liu, Yuanqing Yu, et al. Toolace: Winning the points of llm function calling. *arXiv preprint arXiv:2409.00920*, 2024.

[16] Qiqiang Lin, Muning Wen, Qiuying Peng, Guanyu Nie, Junwei Liao, Jun Wang, Xiaoyun Mo, Jiamu Zhou, Cheng Cheng, Yin Zhao, et al. Hammer: Robust function-calling for on-device language models via function masking. *arXiv preprint arXiv:2410.04587*, 2024.

[17] Jianguo Zhang, Tian Lan, Ming Zhu, Zuxin Liu, Thai Quoc Hoang, Shirley Kokane, Weiran Yao, Juntao Tan, Akshara Prabhakar, Haolin Chen, et al. xlam: A family of large action models to empower ai agent systems. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 11583–11597, 2025.

[18] Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025.

[19] 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. *arXiv preprint arXiv: 2409.19256*, 2024.

[20] Shishir G Patil, Huanzhi Mao, Fanjia Yan, Charlie Cheng-Jie Ji, Vishnu Suresh, Ion Stoica, and Joseph E Gonzalez. The berkeley function calling leaderboard (bfcl): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*.

[21] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.

[22] Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Tianjun Zhang, Erran Li, Raluca Ada Popa, and Ion Stoica. Deepscaler: Surpassing o1-preview with a 1.5b model by scaling rl. https://pretty-radio-b75.notion.site/DeepScaleR-Surpassing-O1-Preview-with-a-1-5B-Model-by-Scaling-RL-19681902c1468005bed8ca303013a4e2, 2025. Notion Blog.

[23] MAA. American invitational mathematics examination - aime. *In American Invitational Mathematics Examination - AIME 2024*, 2024.

[24] MAA. American invitational mathematics examination - amc. *In American Invitational Mathematics Examination - AMC*, 2024.

[25] 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.

[26] 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.

[27] 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.

[28] Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang, Yuchen Zhang, Jiacheng Chen, Wendi Li, Bingxiang He, Yuchen Fan, Tianyu Yu, et al. Process reinforcement through implicit rewards. *arXiv preprint arXiv:2502.01456*, 2025.

[29] Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, et al. Measuring coding challenge competence with apps. *arXiv preprint arXiv:2105.09938*, 2021.

[30] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022.

[31] Rongao Li, Jie Fu, Bo-Wen Zhang, Tao Huang, Zhihong Sun, Chen Lyu, Guang Liu, Zhi Jin, and Ge Li. Taco: Topics in algorithmic code generation dataset. *arXiv preprint arXiv:2312.14852*, 2023.

[32] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Jun-Mei Song, Mingchuan Zhang, Y. K. Li, Yu Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *ArXiv*, abs/2402.03300, 2024.

[33] Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, Jingren Zhou, and Junyang Lin. Group sequence policy optimization. *ArXiv*, abs/2507.18071, 2025.

[34] Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Honglin Yu, Weinan Dai, Yuxuan Song, Xiang Wei, Haodong Zhou, Jingjing Liu, Wei Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yong-Xu Wu, and Mingxuan Wang. Dapo: An open-source llm reinforcement learning system at scale. *ArXiv*, abs/2503.14476, 2025.

[35] Vaishnavi Shrivastava, Ahmed Awadallah, Vidhisha Balachandran, Shivam Garg, Harkirat Singh Behl, and Dimitris Papailiopoulos. Sample more to think less: Group filtered policy optimization for concise reasoning. *ArXiv*, abs/2508.09726, 2025.

[36] Shih-Yang Liu, Xin Dong, Ximing Lu, Shizhe Diao, Mingjie Liu, Min-Hung Chen, Hongxu Yin, Yu-Chiang Frank Wang, Kwang-Ting Cheng, Yejin Choi, Jan Kautz, and Pavlo Molchanov. Dler: Doing length penalty right - incentivizing more intelligence per token via reinforcement learning. *ArXiv*, abs/2510.15110, 2025.

[37] Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. Safe rlhf: Safe reinforcement learning from human feedback. *ArXiv*, abs/2310.12773, 2023.

[38] Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke S. Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. Personalized soups: Personalized large language model alignment via post-hoc parameter merging. *ArXiv*, abs/2310.11564, 2023.

[39] Yuhang Lai, Siyuan Wang, Shujun Liu, Xuanjing Huang, and Zhongyu Wei. Alarm: Align language models via hierarchical rewards modeling. In *Annual Meeting of the Association for Computational Linguistics*, 2024.

[40] DeepSeek-AI, Aixin Liu, Aoxue Mei, Bangcai Lin, Bing Xue, Bing-Li Wang, Bingzheng Xu, Bochao Wu, Bowei Zhang, Chaofan Lin, Chen Dong, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenhao Xu, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Erhang Li, Fangqi Zhou, Fangyun Lin, Fucong Dai, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Hanwei Xu, Hao Li, Haofen Liang, Haoran Wei, Haowei Zhang, Hao sheng Luo, Haozhe Ji, Honghui Ding, Hongxuan Tang, Huanqi Cao, Huazuo Gao, Huixian Qu, Hui Zeng, Jialiang Huang, Jiashi Li, Jiaxin Xu, Jiewen Hu, JingChang Chen, Jingting Xiang, Jingyang Yuan, Jing Cheng, Jinhua Zhu, Jun Ran, Junguang Jiang, Junjie Qiu, Junlong Li, Jun-Mei Song, Kai Dong, Kaige Gao, Kang Guan, Kexin Huang, Kexing Zhou, Kezhao Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Wang, Liang Zhao, Liangsheng Yin, Lihua Guo, Ling-Li Luo, Linwang Ma, Litong Wang, Liyue Zhang, M. S. Di, M. Y. Xu, Mingchuan Zhang, Minghua Zhang, Min Tang, Mingxu Zhou, P. Huang, Peixin Cong, Peiyi Wang, Qiancheng Wang, Qihao Zhu, Qingyang Li, Qinyu Chen, Qiushi Du, Ruiling Xu, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runqiu Yin, Runxin Xu, Ruomeng Shen, Ruoyu Zhang, S. H. Liu, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaofei Cai, Shaoyuan Chen, Shengding Hu, Shengyu Liu, Shiqiang Hu, Shirong Ma, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, Songyang Zhou, Tao Ni, Tao Yun, Tian Pei, Tian Ye, Tianyuan Yue, Wangding Zeng, Wen Liu, Wenfeng Liang, Wenjie Pang, Wenjing Luo, Wenjun Gao, Wentao Zhang, Xi Gao, Xiangwen Wang, Xiaoling Bi, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaokang Zhang, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xingyou Li, Xinyu Yang, Xinyuan Li, Xu Chen, Xuecheng Su, Xuehai Pan, Xuheng Lin, Xuwei Fu, Y. Q. Wang, Yang Zhang, Yanhong Xu, Yanru Ma, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Qian, Yingpu Yu, Yichao Zhang, Yifan Ding, Yifan Shi, Yi Xiong, Ying He, Ying Zhou, Yinmin Zhong, Yishi Piao, Yisong Wang, Yixiao Chen, Yixuan Tan, Yixuan Wei, Yiyang Ma, Yiyuan Liu, Yonglun Yang, Yongqiang Guo, Yongtong Wu, Yu Wu, Yuan Cheng, Yuan Ou, Yuanfan Xu, Yuduan Wang, Yue Gong, Yuhan Wu, Yuheng Zou, Yukun Li, Yunfan Xiong, Yu-Wei Luo, Yu mei You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehua Zhao, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhen guo Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhixian Huang, Zhiyu Wu, Zhuoshu Li, Zhuping Zhang, Zian Xu, Zihao Wang, Zihui Gu, Zijia Zhu, Zi-Rui Li, Zipeng Zhang, Ziwei Xie, Ziyi Gao, Zizheng Pan, Zongqing Yao, Bei Feng, Hui Li, J. L. Cai, Jiaqi Ni, Lei Xu, Meng Li, Ning Tian, R. J. Chen, Ruiqi Jin, S. S. Li, Shuang Zhou, Tianyu Sun, X. Q. Li, Xiangyu Jin, Xiaojin Shen, Xiaosha Chen, Xinnan Song, Xinyi Zhou, Y. X. Zhu, Yanping Huang, Yao Li, Yi Zheng, Yuchen Zhu, Yunxiang Ma, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, Dong-Li Ji, Jian Liang, Jianzhong Guo, Jin Chen, Leyi Xia, Miaojun Wang, Mingming Li, Peng Zhang, Ruyi Chen, Shangmian Sun, Shao-Kang Wu, Shengfeng Ye, T.Wang, W. L. Xiao, Wei An, Xianzu Wang, Xiaowen Sun, Xiaoxiang Wang, Ying Tang, Yukun Zha, Ze-Na Zhang, Zhenghua Ju, Zhen Zhang, and Zihua Qu. Deepseek-v3.2: Pushing the frontier of open large language models. 2025.

[41] Haotian Luo, Li Shen, Haiying He, Yibo Wang, Shiwei Liu, Wei Li, Naiqiang Tan, Xiaochun Cao, and Dacheng Tao. O1-pruner: Length-harmonizing fine-tuning for o1-like reasoning pruning. *ArXiv*, abs/2501.12570, 2025.

[42] Daman Arora and Andrea Zanette. Training language models to reason efficiently. *ArXiv*, abs/2502.04463, 2025.

[43] Jingyang Yi and Jiazheng Wang. Shorterbetter: Guiding reasoning models to find optimal inference length for efficient reasoning. *ArXiv*, abs/2504.21370, 2025.

[44] Pranjal Aggarwal and Sean Welleck. L1: Controlling how long a reasoning model thinks with reinforcement learning. *ArXiv*, abs/2503.04697, 2025.

[45] Jinyan Su and Claire Cardie. Thinking fast and right: Balancing accuracy and reasoning length with adaptive rewards. *ArXiv*, abs/2505.18298, 2025.

# A. Training stability issue of GDPO without batch-wise advantage normalization
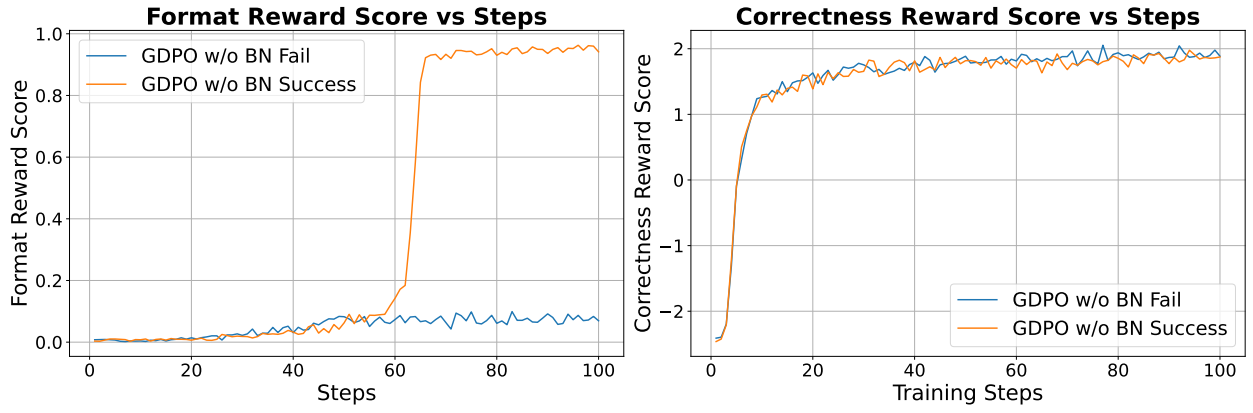


Figure 8 | Training stability of GDPO with and without batch-wise advantage normalization. Runs without normalization occasionally fail to converge.

# B. ToolRL Training Prompt Format

**System Prompt for ToolRL Training**

You are a helpful dialogue assistant capable of leveraging tool calls to solve user tasks and provide structured chat responses.

**Available Tools**
In your response, you can use the following tools:
`{ { Tool List } }`

**Steps for Each Turn**

1. **Think**: Recall relevant context and analyze the current user goal.
2. **Decide on Tool Usage**: If a tool is needed, specify the tool and its parameters.
3. **Respond Appropriately**: If a response is needed, generate one while maintaining consistency across user queries.

**Output Format**
`<think>` Your thoughts and reasoning `</think>`
`<tool_call>` {"name": "Tool name", "parameters": {"Parameter name": "Parameter content", " ... ...": " ... ..."}}
{"name": " ... ...", "parameters": {" ... ...": " ... ...", " ... ...": " ... ..."}}
...
`</tool_call>`
`<response>`AI's final response `</response>`

**Important Notes**

1. You must always include the `<think>` field to outline your reasoning. Provide at least one of `<tool_call>` or `<response>`. Decide whether to use `<tool_call>` (possibly multiple times), `<response>`, or both.
2. You can invoke multiple tool calls simultaneously in the `<tool_call>` fields. Each tool call should be a JSON object with a "name" field and a "parameters" field containing a dictionary of parameters. If no parameters are needed, leave the "parameters" field an empty dictionary.
3. Refer to the previous dialogue records in the history, including the user's queries, previous `<tool_call>`, `<response>`, and any tool feedback noted as `<obs>` (if exists).

**User Prompt for ToolRL Training**

**Dialogue History**

`<user>` { { Initial User Input } } `</user>`

`<think>`  Round 1 Model Thought  `</think>`
{ { Round 1 model output `<tool_call>` or `<response>` } }
`<obs>` Round 1 Observation `</obs>`

... ...

`<user>` { { User Input } } `</user>`

... ...

# C. Tool Calling Reward Functions

**Format Reward.** The format reward $\mathcal{R}_{\text{format}} \in \{0, 1\}$ checks whether the model output satisfies the required structure and contains all necessary fields in the correct order:

$$\mathcal{R}_{\text{format}} = \begin{cases} 1, & \text{if all required fields appear and are in the correct order,} \\ 0, & \text{otherwise.} \end{cases} \tag{9}$$

**Correctness Reward.** The correctness reward $\mathcal{R}_{\text{correct}} \in [-3, 3]$ evaluates the predicted tool calls $P = \{P_1, \ldots, P_m\}$ against the ground-truth calls $G = \{G_1, \ldots, G_n\}$. It consists of three components:

- **Tool Name Matching:**

$$r_{\text{name}} = \frac{|N_G \cap N_P|}{|N_G \cup N_P|} \in [0, 1],$$

  where $N_G$ and $N_P$ are the sets of tool names from ground-truth and predicted calls, respectively.

- **Parameter Name Matching:**

$$r_{\text{param}} = \sum_{G_j \in G} \frac{|\text{keys}(G_j) \cap \text{keys}(P_j)|}{|\text{keys}(G_j) \cup \text{keys}(P_j)|} \in [0, |G|],$$

  where $\text{keys}(G_j)$ and $\text{keys}(P_j)$ are the parameter names of the ground-truth and predicted calls.

- **Parameter Content Matching:**

$$r_{\text{value}} = \sum_{G_j \in G} \sum_{k \in \text{keys}(G_j)} \mathbf{1}[P_G[k] = P_P[k]] \in \left[0, \sum_{G_j \in G} |\text{keys}(G_j)|\right],$$

  where $P_G[k]$ and $P_P[k]$ are the parameter values for the ground-truth and predicted calls.

- **Total Match Score:**

$$r_{\text{match}} = r_{\text{name}} + r_{\text{param}} + r_{\text{value}} \in [0, S_{\text{max}}],$$

  where

$$S_{\text{max}} = 1 + |G| + \sum_{G_j \in G} |\text{keys}(G_j)|.$$

The final correctness reward is computed by finding the optimal matching between $P$ and $G$ to maximize the total match score:

$$\mathcal{R}_{\text{correct}} = 6 \cdot \frac{R_{\text{max}}}{S_{\text{max}}} - 3 \in [-3, 3].$$

where $R_{\text{max}}$ denotes the total match score from the optimal matching.

# D. ToolRL Hyperparameters Setting

Table 6 | GDPO verl training configuration. All hyperparameter settings are kept identical to those used in ToolRL[12].

| Parameter | Value |
|---|---|
| trainer.total_epochs | 15 |
| data.train_batch_size | 512 |
| actor_rollout_ref.actor.ppo_mini_batch_size | 128 |
| data.max_prompt_length | 2048 |
| actor_rollout_ref.actor.optim.lr | 1.00E-06 |
| actor_rollout_ref.rollout.n | 4 |
| algorithm.kl_ctrl.kl_coef | 0.001 |

# E. Math/Coding Reasoning Hyperparameters Setting

Table 7 | GDPO verl training configuration

| Parameter | Value |
|---|---|
| data.train_batch_size | 512 |
| actor_rollout_ref.actor.ppo_mini_batch_size | 64 |
| actor_rollout_ref.actor.ppo_epochs | 1 |
| data.max_prompt_length | 1024 |
| actor_rollout_ref.actor.optim.lr | 1.00E-06 |
| actor_rollout_ref.rollout.temperature | 1 |
| actor_rollout_ref.rollout.n | 16 |
| actor_rollout_ref.actor.clip_ratio_low | 0.2 |
| actor_rollout_ref.actor.clip_ratio_high | 0.28 |
| algorithm.filter_groups.enable | TRUE |
| algorithm.filter_groups.metric | seq_reward |
| actor_rollout_ref.actor.kl_loss_coef | 0.0005 |
| actor_rollout_ref.actor.kl_loss_type | mse |

## F. Training curves of GRPO and GDPO when training DeepSeek-R1-7B and Qwen3-4B-Instruct with $\mathcal{R}_{\text{length}}$ and $\mathcal{R}_{\text{correct}}$ on math reasoning data.



Figure 9 | Training behavior of GRPO and GDPO when optimizing DeepSeek-R1-7B across correctness reward, length reward, and maximum batch response length on math reasoning data. We can see that GDPO maintains improving correctness and better adherence to length constraints over GRPO.



Figure 10 | Training behavior of GRPO and GDPO when optimizing Qwen3-4B-Instruct across correctness reward, length reward, and maximum batch response length on math reasoning data. We can see that GDPO maintains improving correctness and better adherence to length constraints over GRPO.

## G. Comparison of GRPO/GDPO finetuned DeepSeek-R1-7B models under varying length reward weights $\{1.0, 0.75, 0.5, 0.25\}$ with and without the conditioned length reward $\tilde{\mathcal{R}}_{\mathbf{length}}$ on math reasoning tasks

Table 8 | Comparison of GRPO/GDPO finetuned DeepSeek-R1-7B models under varying length reward weights $\{1.0, 0.75, 0.5, 0.25\}$ with the normal length reward $\mathcal{R}_{\mathrm{length}}$ on math reasoning tasks

| | Length Reward Weight | MATH ↑ | Exceed ↓ | AIME ↑ | Exceed ↓ | Amc ↑ | Exceed ↓ | Minerva ↑ | Exceed ↓ | Olympiad ↑ | Exceed ↓ | Avg Acc ↑ | Avg Exceed ↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DeepSeek-R1-7B | - | 93.6% | 26.0% | 55.4% | 85.6% | 82.9% | 57.2% | 49.8% | 41.8% | 58.2% | 60.6% | 68.0% | 54.2% |
| GRPO-$\mathcal{R}_{\mathrm{length}}$ | 0.25 | 94.0% | 0.5% | 53.3% | 4.9% | 83.8% | 0.6% | 53.2% | 1.9% | 59.8% | 2.4% | 68.8% | 2.1% |
| | 0.5 | 94.2% | 0.6% | 52.1% | 2.3% | 83.2% | 0.8% | 53.9% | 0.2% | 60.2% | 0.8% | 68.7% | 0.9% |
| | 0.75 | 93.5% | 0.3% | 51.2% | 2.5% | 83.0% | 0.5% | 53.4% | 0.1% | 58.2% | 1.5% | 67.8% | 1.0% |
| | 1.0 | 94.1% | 0.5% | 50.2% | 2.1% | 83.8% | 0.6% | 53.2% | 0.2% | 60.2% | 1.1% | 68.3% | 0.9% |
| GDPO-$\mathcal{R}_{\mathrm{length}}$ | 0.25 | 94.2% | 0.5% | 54.7% | 3.9% | 84.5% | 1.4% | 54.1% | 1.1% | 59.1% | 1.3% | 69.3% | 1.6% |
| | 0.5 | 93.2% | 0.2% | 53.8% | 0.8% | 84.1% | 0.4% | 53.3% | 0.2% | 58.5% | 1.1% | 68.6% | 0.5% |
| | 0.75 | 93.6% | 0.2% | 54.4% | 1.5% | 83.4% | 0.2% | 53.4% | 0.3% | 58.8% | 0.6% | 68.7% | 0.5% |
| | 1.0 | 93.9% | 0.1% | 53.1% | 0.2% | 84.0% | 0.3% | 53.8% | 0.1% | 59.3% | 0.4% | 68.8% | 0.2% |

Table 9 | Comparison of GRPO/GDPO finetuned DeepSeek-R1-7B models under varying length reward weights $\{1.0, 0.75, 0.5, 0.25\}$ with the conditioned length reward $\tilde{\mathcal{R}}_{\mathrm{length}}$ on math reasoning tasks

| | Length Reward Weight | MATH ↑ | Exceed ↓ | AIME ↑ | Exceed ↓ | Amc ↑ | Exceed ↓ | Minerva ↑ | Exceed ↓ | Olympiad ↑ | Exceed ↓ | Avg Acc ↑ | Avg Exceed ↓ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DeepSeek-R1-7B | - | 93.6% | 26.0% | 55.4% | 85.6% | 82.9% | 57.2% | 49.8% | 41.8% | 58.2% | 60.6% | 68.0% | 54.2% |
| GRPO-$\tilde{\mathcal{R}}_{\mathrm{length}}$ | 0.25 | 94.4% | 6.4% | 56.8% | 51.2% | 85.8% | 18.2% | 54.5% | 6.5% | 61.2% | 18.8% | 70.6% | 20.2% |
| | 0.5 | 94.2% | 4.0% | 56.5% | 43.8% | 85.9% | 13.6% | 54.8% | 7.1% | 60.6% | 17.2% | 70.4% | 17.1% |
| | 0.75 | 94.1% | 3.4% | 55.5% | 38.3% | 85.2% | 11.2% | 54.3% | 5.7% | 60.8% | 14.3% | 70.0% | 14.6% |
| | 1.0 | 93.2% | 2.7% | 53.3% | 29.2% | 82.9% | 8.6% | 53.2% | 2.5% | 59.1% | 10.3% | 68.3% | 10.7% |
| GDPO-$\tilde{\mathcal{R}}_{\mathrm{length}}$ | 0.25 | 94.9% | 5.1% | 57.7% | 32.8% | 86.2% | 11.2% | 54.9% | 5.1% | 62.1% | 14.6% | 71.2% | 13.8% |
| | 0.5 | 94.5% | 3.1% | 57.7% | 32.1% | 85.8% | 9.3% | 54.1% | 3.9% | 61.6% | 13.3% | 70.8% | 12.3% |
| | 0.75 | 94.6% | 2.8% | 56.0% | 31.9% | 86.9% | 10.3% | 53.0% | 2.6% | 61.3% | 12.8% | 70.4% | 12.1% |
| | 1.0 | 93.9% | 1.4% | 57.7% | 12.3% | 85.9% | 3.8% | 53.4% | 1.0% | 60.8% | 6.2% | 70.4% | 4.9% |