# Rectify Evaluation Preference: Improving LLMs' Critique on Math Reasoning via Perplexity-aware Reinforcement Learning

**Changyuan Tian**[1,2,3,*], **Zhicong Lu**[1,2,3*], **Shuang Qian**[4], **Nayu Liu**[5], **Peiguang Li**[4,†],
**Li Jin**[1,†], **Leiyi Hu**[1,2,3], **Zhizhao Zeng**[4], **Sirui Wang**[4], **Ke Zeng**[4], **Zhi Guo**[1]

[1]Aerospace Information Research Institute, Chinese Academy of Sciences
[2]Key Laboratory of Target Cognition and Application Technology (TCAT)
[3]School of Electronic, Electrical and Communication Engineering, University of Chinese Academy of Sciences
[4]Meituan
[5]School of Computer Science and Technology, Tiangong University
tianchangyuan21@mails.ucas.edu.cn, jinlimails@gmail.com

## Abstract

To improve Multi-step Mathematical Reasoning (MsMR) of Large Language Models (LLMs), it is crucial to obtain scalable supervision from the corpus by automatically critiquing mistakes in the reasoning process of MsMR and rendering a final verdict of the problem-solution. Most existing methods rely on crafting high-quality supervised fine-tuning demonstrations for critiquing capability enhancement and pay little attention to delving into the underlying reason for the poor critiquing performance of LLMs. In this paper, we orthogonally quantify and investigate the potential reason — imbalanced evaluation preference, and conduct a statistical preference analysis. Motivated by the analysis of the reason, a novel perplexity-aware reinforcement learning algorithm is proposed to rectify the evaluation preference, elevating the critiquing capability. Specifically, to probe into LLMs' critiquing characteristics, a One-to-many Problem-Solution (OPS) benchmark is meticulously constructed to quantify the behavior difference of LLMs when evaluating the problem solutions generated by itself and others. Then, to investigate the behavior difference in depth, we conduct a statistical preference analysis oriented on perplexity and find an intriguing phenomenon — "LLMs incline to judge solutions with lower perplexity as correct", which is dubbed as imbalanced evaluation preference. To rectify this preference, we regard perplexity as the baton in the algorithm of Group Relative Policy Optimization, supporting the LLMs to explore trajectories that judge lower perplexity as wrong and higher perplexity as correct. Extensive experimental results on our built OPS and existing available critic benchmarks demonstrate the validity of our method.

## Introduction

Large Language Models (LLMs) have demonstrated exceptional capabilities in handling a wide range of tasks (Grattafiori et al. 2024; Yang et al. 2025; Lu et al. 2025; Liu et al. 2025a; Lu et al. 2023; Jia et al. 2025). However, their performance in Multi-step Mathematical Reasoning

(MsMR) remains relatively limited. To improve the MsMR of LLMs, it is crucial to obtain scalable supervision from the corpus. Given the labor-intensive of directly annotating MsMR, researchers focus on employing instructed LLMs to automatically critiquing mistakes in the reasoning process of MsMR and rending a final verdict of the problem-solution.

However, the critiquing performance of instructed LLMs is far from meeting actual needs. To address this issue, existing works focus on crafting high-quality supervised fine-tuning demonstrations for critiquing capability enhancement. For instance, Wang et al. (2024a) adopt the monte carlo sampling to scalarly label the sample and achieve a critiquing model, which only judges the correctness of the entire problem-solution without critiquing process. To improve the interpretability of judgment, Gao et al. (2025) leverages GPT-4 to finely label the intermediate step of sample, encouraging the critiquing model to identify the mistakes in the reasoning process. However, these data-driven methods pay little attention to delving into the underlying reason for the poor critiquing performance of instructed LLMs.

To bridge the gap of existing methods, our work is initiated with probing into LLMs' critiquing characteristics. Concretely, a One-to-many Problem Solution (OPS) benchmark is meticulously constructed, where each sample comprises one mathematical problem and many solutions generated by various families of LLMs (e.g., LLaMA, Qwen, and Mistral). Later, we quantify the behavior difference of LLMs when evaluating the problem-solution generated by itself and others. To investigate the behavior difference in depth, we conduct a statistical preference analysis oriented on perplexity and find an intriguing phenomenon — "LLMs incline to judge solutions with lower perplexity as correct", which is dubbed as the imbalanced evaluation preference. This biased preference hinders the model mining true cause-effect from supervised demonstrations during training (e.g., whether the final verdict of "correct" comes from lower perplexity or problem-solution itself.)

Motivated by the analysis of the potential reasons (i.e., the imbalanced evaluation preference) for the poor critiquing performance of LLMs, in this paper, we propose

a novel perplexity-aware reinforcement learning algorithm (i.e., Group Relative Policy Optimization, GRPO) to rectify the evaluation preference. Specifically, we leverage perplexity as the baton to rescale the advantage distribution during Reinforcement Learning (RL), assigning greater weight to counter-preference trajectories (e.g., low perplexity but predicted wrong), thereby balancing LLMs' exploration. Additionally, we apply class-level loss aggregation by independently aggregating the losses for wrong-label and correct-label trajectories, ensuring that the rescaled comparative advantages within each class are faithfully reflected during optimization. These two designs jointly rectify evaluation preference and consequently improve the RL fine-tuning performance.

To validate the effectiveness of our method, we carry out extensive experiments on our built OPS and existing available critic benchmarks with mainstream LLMs. Compared to standard GRPO, our method alleviates the imbalanced evaluation preference and exhibits better critiquing capability. In summary, our main contributions include:

- We meticulously construct a One-to-many Problem Solution (OPS) benchmark to quantify and investigate the potential reason (i.e., imbalanced evaluation preference) for the poor critiquing performance of LLMs.
- Motivated by the analysis ("LLMs incline to judge solutions with lower perplexity as correct") of the potential reasons, we propose a novel perplexity-aware reinforcement learning algorithm to rectify the imbalanced evaluation preference, which supports the LLMs to explore counter-preference trajectories that judge lower perplexity as wrong and higher perplexity as correct.
- Extensive experiments on our built OPS benchmark demonstrate our method effectively alleviates the evaluation preference. The newly achieved state-of-the-art performance on existing available critic benchmarks further demonstrates the validity of our method.

## Analysis of Evaluation Preference

In this section, we delve into the underlying reason for the poor critiquing performance of LLMs. To probe into the critiquing characteristics of LLMs, we meticulously construct a **O**ne-to-many **P**roblem-**S**olution (OPS) benchmark, to quantify the behavioral differences when LLMs evaluate problem solutions generated by themselves and by others. Furthermore, we conduct a statistical preference analysis based on perplexity and uncover an intriguing phenomenon: LLMs tend to judge solutions with lower perplexity as correct, a bias we refer to as *imbalanced evaluation preference*.

### One-to-Many Problem-Solution Benchmark

Unlike previous benchmarks that construct evaluation sets by randomly collecting mathematical reasoning solutions, our OPS benchmark specifically focuses on solutions generated by different families of LLMs that solve the same problem and arrive at the same final answer. This controlled design ensures that observed behavioral differences can be attributed to variations in reasoning processes.

**Mathematical Problems and Solutions Collection**. The mathematical problems used in OPS benchmark are sourced from the MATH (Hendrycks et al. 2021) test dataset, which features a wide range of problem types and difficulty. To ensure diversity in the generated solutions, we employ three distinct LLMs—Qwen2-7B-Instruct, LLaMA3.1-8B-Instruct, and Mistral-7B-Instruct-v0.3—to independently solve each problem, using a chain-of-thought prompting strategy with a temperature setting of 0.7. Solution correctness is labeled by comparing the final answer to the gold answer provided in the MATH dataset; a solution is labeled as *correct* if the answers match, and *wrong* otherwise. To further ensure data quality, we additionally leverage an advanced discriminative process reward model, Qwen2.5-Math-PRM-72B (Zhang et al. 2025b), to filter the generated solutions. Solutions for which the correctness labels assigned by the process reward model and by answer matching are inconsistent are removed from the dataset.

**Construction of Solution Triples**. To ensure that observed behavioral differences can be attributed to variations in reasoning processes, we construct solution triples. Each triple consists of independently generated solutions by Qwen, LLaMA, and Mistral for the same problem and yielding the same final answer. Formally, each triple is defined as:

$$T = [x, \{y_1, y_2, y_3\}, a, v] \tag{1}$$

where $x$ denotes the mathematical problem, $y_i$ is the solution generated by the $i$-th model ($i \in \{1, 2, 3\}$ corresponding to Qwen, LLaMA, and Mistral), $a$ is the final answer extracted from $y_i$, $v \in \{\text{correct}, \text{wrong}\}$ indicates the correctness of $y_i$. To achieve a balanced evaluation, we ensure that the dataset contains an equal number of 'correct'-label and 'wrong'-label solutions (i.e., a 1:1 ratio).

Consequently, our OPS benchmark consists of three test subsets—LLaMA3.1-8B-Instruct, Mistral-7B-Instruct-v0.3, and Qwen2-7B-Instruct—each containing 630 samples with the same problems and final answers, resulting in a total of 1,890 test samples.

**Metrics**. To measure evaluation preference, we introduce the *Balance Indicator* (BI), defined as the difference between the false positive rate (FPR) and the false negative rate (FNR):

$$\text{BI} = \text{FPR} - \text{FNR} \tag{2}$$

here, FPR is the proportion of actually incorrect solutions erroneously judged as correct, while FNR is the proportion of actually correct solutions erroneously judged as wrong. A BI close to zero indicates balanced evaluation preference; positive BI suggests overestimation of incorrect solutions, and negative BI suggests underestimation of correct ones. The BI ranges from $-1$ (maximal underestimation) to $1$ (maximal overestimation), with $0$ indicating perfect balance. In addition, we also use basic accuracy as another metric.

### Imbalanced Evaluation Preference

Based on OPS benchmark, we evaluate Qwen2-7B-Instruct, LLaMA3.1-8B-Instruct, and Mistral-7B-Instruct-v0.3, with

| Model | L-subset | | M-subset | | Q-subset | |
|---|---|---|---|---|---|---|
| | Acc | BI | Acc | BI | Acc | BI |
| Qwen2-7B-Instruct | 65.56 | 34.60 | 67.94 | 2.54 | 60.63 | 66.03 |
| LLaMA3.1-8B-Instruct | 64.13 | 17.78 | 62.86 | -30.48 | 65.24 | 16.83 |
| Mistral-7B-Instruct-v0.3 | 52.70 | 81.27 | 51.90 | 86.03 | 52.54 | 86.67 |

Table 1: Performance of each model on self and cross-evaluation per subset. The L-subset, M-subset, and Q-subset columns correspond to evaluation results on the LLaMA3.1-8B-Instruct, Mistral-7B-Instruct-v0.3, and Qwen2-7B-Instruct test subsets, respectively.

the results presented in Table 1. Firstly, a direct observation is that despite the three subsets sharing the same problems and final answers, the three LLMs exhibit markedly different evaluation performance across the subsets. This suggests that evaluation behavior may be affected by model-specific reasoning, leading to inconsistent assessments across solutions generated by different models. Furthermore, we observe a clear imbalanced evaluation preference in LLMs: each model tends to achieve a higher and more positive BI when evaluating its own solutions, while BI drops substantially when evaluating solutions from other models (e.g., Qwen: self 66.03 vs. cross 2.54, LLaMA: self 17.78 vs. cross $-30.48$, Mistral: self 86.03 vs. cross 81.27). This imbalanced preference impairs model evaluation performance.

### Statistical Preference Analysis

Motivated by the observed evaluation preference, we further conduct a fine-grained statistical analysis by exploring the correlation between perplexity and BI. Perplexity reflects how closely a solution aligns with the critic model's own generation style, with lower values indicating smaller textual divergence. Specifically, for each critic model, we calculate the perplexity of each problem-solution pair in OPS benchmark, partition the data into decile bins based on perplexity, and perform linear regression analysis to investigate the relationship between perplexity and BI.
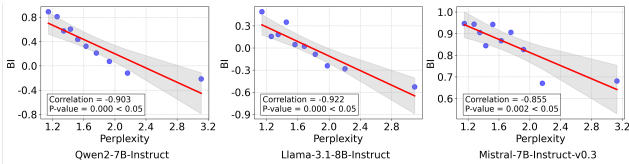


Figure 1: Visualization of the negative correlation between perplexity and BI for each critic model. Each point represents a decile bin, with the red regression line and gray-shaded area denoting the fitted linear trend and the 95% confidence interval, respectively. Higher perplexity values are associated with lower BI.

We observe a clear linear negative correlation between perplexity and BI (see Figure 1): as perplexity increases—indicating that a solution deviates more from the critic model's generation style—the critic model is more likely to judge it as wrong, resulting in BI values shifting

toward $-1$. This pattern holds consistently across all three models, each exhibiting a significant negative correlation coefficient with $p$-values less than 0.05. The strong correlation between evaluation preference and perplexity suggests a promising entry point for us to manipulate and mitigate such bias.

## Method

### Problem Setup

Given a mathematical problem $x$ and its proposed solution $y$, our model is designed to provide a comprehensive critique of $y$. Specifically, the model outputs a triplet $(f, v, \hat{e})$, where $f$ denotes the step-by-step critique on each reasoning step in $y$, $v \in \{\text{correct, wrong}\}$ indicates the overall correctness of the solution, and $\hat{e}$ records the content of the first erroneous step, if any. If all steps in $y$ are correct, then $v = \text{correct}$ and $\hat{e} = \varnothing$; otherwise, $v = \text{wrong}$ and $\hat{e}$ contains the first incorrect step identified by the model. Formally, this evaluation process can be expressed as:

$$(f, v, \hat{e}) = LLM(x, y) \tag{3}$$

where $LLM(\cdot)$ denotes the evaluation model. This formulation enables fine-grained analysis of the solution process, supporting both detailed diagnostic feedback and precise error localization.

### Training Data Curation

The mathematical problems used in our study are sourced from the MATH (Hendrycks et al. 2021) train dataset, which covers diverse problem types and difficulty levels. Consistent with the OPS benchmark, we employ the same three LLMs—Qwen2-7B-Instruct, LLaMA3.1-8B-Instruct, and Mistral-7B-Instruct-v0.3—for solution sampling, ensuring a wide range of solution perplexities.

A key requirement for our training data curation is the accurate identification of the first erroneous step within each solution. Following the ProcessBenchmark methodology (Zheng et al. 2025), we annotate each solution as follows: first, Qwen2.5-72B-Instruct (Qwen et al. 2025) is used to segment the solution into individual reasoning steps. Each step is then evaluated by the state-of-the-art process reward model, Qwen2.5-Math-PRM-72B (Zhang et al. 2025b). Solutions that exhibit inconsistencies between stepwise and overall correctness are filtered out to maintain annotation quality. The first step with a reward score below 0.8 is labeled as the initial error.

Formally, each training example is represented as $D = \{x, y, a, v, e^*\}$, where $x$ is the mathematical problem, $y$ is the solution generated by the model (Qwen, LLaMA, or Mistral), $a$ is the gold answer, $v \in \{\text{correct, wrong}\}$ denotes the correctness of $y$, and $e^*$ is the first erroneous step in $y$.

The final training set comprises 5,760 samples, with equal numbers of solutions sampled from each of the three models. For each model, correct and incorrect solutions are balanced, ensuring a well-balanced training dataset.
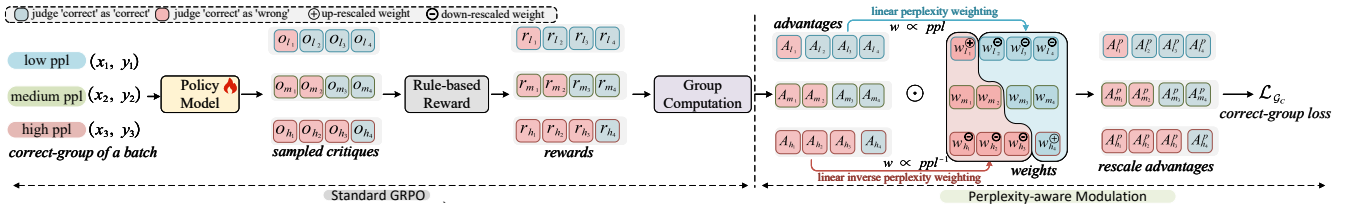
Figure 2: Illustration of our proposed perplexity-aware GRPO. Here, we take the correct-group within a batch (i.e., the sample set where the ground truth is correct) as an example; the wrong group follows the same process.

## Perplexity-aware Reinforcement Learning

To rectify evaluation preference, our approach leverages perplexity to modulate the advantage distribution during reinforcement learning, thereby achieving more balanced evaluation performance. Without loss of generality, we adopt the well-established GRPO (Shao et al. 2024) method as our base, which does not require step-by-step supervision signals. An illustration of our proposed perplexity-aware GRPO is shown in Figure 2.

**Standard GRPO.** GRPO first samples a group of candidate critiques $\{o_i\}_{i=1}^G$ for each question-answer pair, assigns reward scores $r_i$ to these critiques, and then estimates their advantages $A_i$ by normalizing the group-level rewards. Its objective function for each critique is defined as:

$$\mathcal{J}(\theta) = \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left[ \min\left( r_{i,t}(\theta) A_{i,t}, \ \mathrm{clip}\big(r_{i,t}(\theta), \right. \right.$$
$$\left. \left. 1-\varepsilon, 1+\varepsilon\big) A_{i,t} \right) - \beta \, D_{\mathrm{KL}}\big(\pi_\theta \| \pi_{\mathrm{ref}}\big) \right] \tag{4}$$

here, $r_{i,t}(\theta)$ denotes the ratio of the current policy's probability to the old policy's probability for the $t$-th token in the $i$-th critique. The advantage $A_i$ is computed by normalizing the $i$-th critique's reward $r_i$ against the mean and standard deviation of rewards across all $G$ critiques in the group:
$A_i = \frac{r_i - \mathrm{mean}(\{r_i\}_{i=1}^G)}{\mathrm{std}(\{r_i\}_{i=1}^G)}$.

**Perplexity-aware Modulation.** The prior analysis reveals that LLMs tend to evaluate solutions with lower perplexity as correct, and those with higher perplexity as incorrect. This imbalanced evaluation preference narrows the LLMs' potential trajectory exploration space: they rarely explore trajectories that identify high-perplexity solutions as correct or low-perplexity ones as wrong. Such insufficient and imbalanced exploration ultimately leads to suboptimal performance.

To promote more balanced exploration, we introduce a perplexity-aware modulation mechanism. The core idea is to upweight high-perplexity-but-predicted-correct and low-perplexity-but-predicted-wrong cases, thereby encouraging the model to explore counter-preference behaviors in both directions. First, each training batch is split into two groups based on the ground-truth labels: the correct-group $\mathcal{G}_C$, containing cases labeled as correct, and the wrong-group $\mathcal{G}_W$,

containing cases labeled as wrong. Within each group (e.g., $\mathcal{G}_C$), we further partition the data into subgroups according to the model's predictions: $\mathcal{G}_C^c$ for cases predicted as correct, and $\mathcal{G}_C^w$ for cases predicted as wrong. For subgroups $\mathcal{G}_C^c$ and $\mathcal{G}_W^c$, we employ a linear, perplexity-aware advantage modulation, assigning greater weights (i.e., $> 1$) to samples with higher perplexity. Conversely, for subgroups $\mathcal{G}_C^w$ and $\mathcal{G}_W^w$, we apply an inverse linear modulation, whereby lower perplexity corresponds to greater weights.

Formally, take $\mathcal{G}_C$ as an example. Let $ppl$ denote the perplexity of a problem-solution pair within this group. The weight assigned to each corresponding critique is given by

$$w_i = \begin{cases} ppl_i/\mathrm{mean}(\{ppl_k \mid k \in \mathcal{G}_C^c\}), & i \in \mathcal{G}_C^c \\[2mm] \mathrm{mean}(\{ppl_k \mid k \in \mathcal{G}_C^w\})/ppl_i, & i \in \mathcal{G}_C^w \end{cases} \tag{5}$$

where $\mathcal{G}_C^c$ and $\mathcal{G}_C^w$ denote the subgroups of $\mathcal{G}_C$ predicted as correct and wrong, respectively.

These modulation weights $w_i$ are then used to rescale the advantages for each critique prior to policy optimization:

$$A_i^{\mathrm{p}} = w_i \cdot A_i. \tag{6}$$

this perplexity-aware modulation encourages the policy to cover less-explored trajectory spaces, mitigating the LLM's bias and thereby enhancing RL fine-tuning performance.

**Class-Level Loss Aggregation**. In vanilla GRPO, the loss for each sample in a batch is aggregated by taking the overall mean, i.e., $\mathcal{L} = \frac{1}{B} \sum_{i=1}^B \mathcal{J}_i$. In our method, however, the advantage distribution within each training batch is modulated separately within groups defined by the ground truth: the correct group ($\mathcal{G}_C$) and the wrong group ($\mathcal{G}_W$). To ensure that the perplexity-modulated comparative advantages are fully reflected in the loss, we apply class-level loss aggregation. Specifically, we first compute the mean of nonzero loss terms within each group, and then average these group-wise means to obtain the final loss. Formally, the class-level loss is defined as:

$$\mathcal{L}' = \frac{1}{2} \left( \frac{1}{|\mathcal{G}_C'|} \sum_{i \in \mathcal{G}_C'} \mathcal{J}_i + \frac{1}{|\mathcal{G}_W'|} \sum_{i \in \mathcal{G}_W'} \mathcal{J}_i \right), \tag{7}$$

where $\mathcal{G}_C' = \{i \in \mathcal{G}_C \mid \mathcal{J}_i \neq 0\}$ and $\mathcal{G}_W' = \{i \in \mathcal{G}_W \mid \mathcal{J}_i \neq 0\}$ denote the sets of samples with nonzero loss in the correct and wrong groups, respectively.

**Reward Design.** We design a rule-based reward function to guide the model toward generating accurate verdicts and identifying the first erroneous step. The overall reward comprises two components: the format reward $r_f$ and the answer reward $r_a$. The format reward assigns $0.1$ if the response strictly follows the predefined format, and $0$ otherwise.

The answer reward $r_a$ reflects the hierarchical structure of the mathematical reasoning evaluation task. For actually incorrect problem-solution pairs ($v^* =$ 'wrong'), partial reward is given for correctly judging the solution as wrong, with additional reward for accurately localizing the first erroneous step. For actually correct pairs ($v^* =$ 'correct'), the full reward is granted only for correctly judging the solution as correct. Specifically, the answer reward $r_a$ is defined as:

$$
r_a = \begin{cases} 0.8, & v^* = \text{'correct' and } v = \text{'correct'} \\ 0.6 + \text{F1}(\hat{e}, e^*), & v^* = \text{'wrong' and } v = \text{'wrong'} \\ 0, & \text{otherwise} \end{cases}
$$
(8)

where $\text{F1}(\hat{e}, e^*)$ denotes the F1 score between the predicted and ground truth first erroneous steps. To enhance reward stability and eliminate the influence of minor fluctuations, we round the F1 score to one decimal place, thereby suppressing variations smaller than $0.1$. Note that directly employing a simple binary (0-1) reward incentivizes the model to overpredict 'correct' due to the relative ease of obtaining full reward for correct solutions, while failing to promote precise identification of the first erroneous step. To avoid such shortcut behavior, the reward function is designed with a graded scheme (with base rewards such as 0.8 and 0.6 empirically chosen within a typical range) that aligns reward magnitude with task difficulty. The final reward $r$ is computed as $r_f + r_a$.

## Experiments

### Datasets

In our experiments, we employ two benchmarks: our proposed OPS benchmark and the public ProcessBenchmark (Zheng et al. 2025). The OPS benchmark comprises solutions to MATH test data problems generated by Qwen2-7B-Instruct (Team 2024), LLaMA3.1-8B-Instruct (Grattafiori et al. 2024), and Mistral-v0.3-Instruct (Jiang et al. 2023). All subsets share the same problems and final answers, but differ in their solution processes. ProcessBench assesses a model's ability to localize errors in mathematical reasoning and consists of four subsets—GSM8K, MATH, Olympiad-Bench, and Omni-MATH—with solution processes sampled from various models.

### Baselines

We employ three mainstream open-source aligned models: Qwen2-7B-Instruct, LLaMA3.1-8B-Instruct, and Mistral-7B-Instruct-v0.3, to validate the effectiveness of our proposed method. Furthermore, we incorporate SFT, standard GRPO (Shao et al. 2024), and DrGRPO (Liu et al. 2025b) variants as baseline models to facilitate comparison with our perplexity-aware GRPO approach. For SFT, step-by-step supervision is provided by annotations from GPT-4o.

### Experimental Settings

The implementations of the SFT, GRPO, and DrGRPO baselines are based on the publicly available VeRL (Sheng et al. 2025) open-source project. Likewise, our proposed perplexity-aware GRPO method is also developed on top of VeRL. All experiments are conducted using eight NVIDIA A100-80G GPUs with bfloat16 precision. For the training of GRPO, DrGRPO, and Perplexity-aware GRPO, we set the learning rate to $1e-6$ and use a train batch size of 128. The total number of training steps is set to 400. During the rollout phase, 5 samples are generated for each prompt with a temperature of 1.0. Rollout generation is support by vLLM (Kwon et al. 2023). We save checkpoints every 50 steps during training and select the checkpoint with the best average performance across the two benchmarks. For SFT, we train for 3 epochs with a learning rate of $2e-6$ and select the checkpoint with the best average performance. All evaluations use greedy decoding (temperature=0).

### Experimental Results

The experimental results on our proposed OPS benchmark are presented in Table 2. Our perplexity-aware GRPO method outperforms the baseline approaches in terms of both accuracy and balance indicator in most cases across all three base models. These results demonstrate the effectiveness of our method in rectifying evaluation preference and thereby improving RL fine-tuning performance, as well as its generalizability across different model series.

Further, we observe that: (1) SFT demonstrates robust and consistent improvements (e.g., Qwen: 69.00, LLaMA: 71.48, Mistral: 68.31) across models with varying initial performance, whereas RL methods such as GRPO are more sensitive to the base model's starting accuracy and do not offer clear advantages over SFT when the initial performance is low (e.g., Mistral-7B-Instruct-v0.3: starting accuracy 52.38, close to random guess; after RL, 68.41, similar to SFT's 68.31). This suggests that the initial behavior of the base model, such as its evaluation preference, may affect the effectiveness of subsequent RL fine-tuning. (2) For strong base models such as Qwen2 and LLaMA3.1, the GRPO and DrGRPO methods achieve remarkable performance gains over SFT by leveraging the exploration-exploitation capabilities of RL; however, they face severe imbalance issue—larger |BI| and noticeable performance differences across subsets—resulting in suboptimal performance. (3) In contrast, our perplexity-aware GRPO delivers more balanced and superior evaluation results, with smaller |BI| and reduced performance differences across subsets. This highlights the effectiveness of using perplexity to calibrate imbalanced advantage distributions.

The experimental results on the public ProcessBench dataset are summarized in Table 3, where the F1 score is reported based on the accuracies of label-wrong and label-correct samples. Considering ProcessBench's requirement to precisely localize the first erroneous step, our proposed perplexity-aware GRPO achieves superior performance in most cases (11 out of 15) across three base models. These results further demonstrate that our method sub-

| Model | LLaMA3.1-8B-Instruct-Subset | | | | Mistral-7B-Instruct-v0.3-Subset | | | | Qwen2-7B-Instruct-Subset | | | | Average | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | FPR | FNR | BI | Acc | FPR | FNR | BI | Acc | FPR | FNR | BI | Acc | $|BI| \downarrow$ | $Acc \uparrow$ |
| Qwen2-7B-Instruct | 51.75 | 17.14 | 34.60 | 65.56 | 33.33 | 30.79 | 2.54 | 67.94 | 72.38 | 6.35 | 66.03 | 60.63 | 34.39 | 64.71 |
| + SFT | 31.11 | 26.98 | 4.13 | 70.95 | 16.83 | 49.84 | -33.02 | 66.67 | 38.73 | 22.54 | 16.19 | 69.37 | 17.78 | 69.00 |
| + Vanilla GRPO | 7.94 | 39.05 | -31.11 | 76.51 | 3.17 | 54.92 | -51.75 | 70.95 | 11.75 | 34.60 | -22.86 | 76.83 | 35.24 | 74.76 |
| + DrGRPO | 9.52 | 33.33 | -23.81 | 78.57 | 3.81 | 48.25 | -44.44 | 73.97 | 10.79 | 31.43 | -20.63 | 78.89 | 29.63 | 77.14 |
| + Perplexity-aware GRPO | 14.29 | 27.62 | -13.33 | 79.05 | 5.71 | 33.97 | -28.25 | 80.16 | 14.92 | 24.76 | -9.84 | 80.16 | 17.14 | 79.79 |
| LLaMA3.1-8B-Instruct | 44.76 | 26.98 | 17.78 | 64.13 | 21.90 | 52.38 | -30.48 | 62.86 | 43.17 | 26.35 | 16.83 | 65.24 | 21.70 | 64.07 |
| + SFT | 31.11 | 25.71 | 5.40 | 71.59 | 9.52 | 49.21 | -39.68 | 70.63 | 31.75 | 23.81 | 7.94 | 72.22 | 17.67 | 71.48 |
| + Vanilla GRPO | 6.98 | 39.68 | -32.70 | 76.67 | 5.08 | 48.25 | -43.17 | 73.33 | 6.35 | 40.32 | -33.97 | 76.67 | 36.61 | 75.56 |
| + DrGRPO | 6.98 | 38.10 | -31.11 | 77.46 | 6.98 | 44.44 | -37.46 | 74.29 | 7.94 | 35.24 | -27.30 | 78.41 | 31.96 | 76.72 |
| + Perplexity-aware GRPO | 18.10 | 26.67 | -8.57 | 77.62 | 8.57 | 37.46 | -28.89 | 76.98 | 14.29 | 27.30 | -13.02 | 79.21 | 16.83 | 77.94 |
| Mistral-7B-Instruct-v0.3 | 87.94 | 6.67 | 81.27 | 52.70 | 91.11 | 5.08 | 86.03 | 51.90 | 90.79 | 4.13 | 86.67 | 52.54 | 84.65 | 52.38 |
| + SFT | 28.25 | 31.43 | -3.17 | 70.16 | 20.32 | 43.49 | -23.17 | 68.10 | 40.63 | 26.03 | 14.60 | 66.67 | 13.65 | 68.31 |
| + Vanilla GRPO | 14.29 | 44.44 | -30.16 | 70.63 | 13.65 | 53.65 | -40.00 | 66.35 | 17.78 | 45.71 | -27.94 | 68.25 | 32.70 | 68.41 |
| + DrGRPO | 13.02 | 49.21 | -36.19 | 68.89 | 11.43 | 61.59 | -50.16 | 63.49 | 13.97 | 49.52 | -35.56 | 68.25 | 40.64 | 66.88 |
| + Perplexity-aware GRPO | 16.51 | 44.13 | -27.62 | 69.68 | 17.78 | 47.94 | -30.16 | 67.14 | 22.22 | 38.41 | -16.19 | 69.68 | 24.66 | 68.83 |

Table 2: Performance comparison on the OPS benchmark. FPR: false positive rate; FNR: false negative rate; BI: signed balance indicator; |BI|: absolute BI. Best results are in bold; second-best are underlined.

stantially enhances critiquing capability, yielding consistent improvements across four datasets of varying difficulty. Notably, the SFT method lags significantly behind other RL methods, indicating poor error step localization capability, which may be related to the challenge of obtaining large-scale, accurately labeled step-by-step supervision data. In contrast, vanilla GRPO and DrGRPO tend to overfit the training data stem from MATH, as reflected by their high performance on MATH (e.g., LLaMA3.1-8B-Instruct + GRPO achieves 43.38), but they struggle to generalize to other datasets, resulting in suboptimal overall performance. By rectifying imbalanced preference, our perplexity-aware GRPO demonstrates more balanced cross-dataset generalization and achieves the best overall performance.

**Ablation Study**

To further evaluate the effectiveness of our method, we conduct ablation studies on the OPS benchmark using Qwen2-7B-Instruct as the base model. As shown in Table 4, we consider the following three ablation variants: (1) *-w/o class-level aggregation*, which disables class-level aggregation and reverts to vanilla aggregation; (2) *-w/o perplexity modulation*, which removes perplexity-aware modulation; and (3) *-w/o both*, which disables both components and is equivalent to the vanilla GRPO. We make the following observations: (1) removing class-level aggregation leads to an increase in |BI| and a slight decrease in accuracy, indicating that class-level aggregation plays a positive role in maintaining optimization balance. (2) in contrast, removing perplexity modulation results in a significant drop in accuracy, even though it noticeably reduces |BI|. This suggests that simply minimizing |BI| through class-level aggregation does not directly result in substantial performance improvement, as it fails to tackle the core issue of imbalanced exploitation during RL. Fundamentally, leveraging perplexity to balance the model's exploration is crucial for improving RL performance. (3)

disabling both components leads to the worst performance, underscoring the necessity of each core design.

| Model | GM | MT | OB | OM | AVG |
|---|---|---|---|---|---|
| Qwen2-7B-Instruct | 24.97 | 31.47 | 27.36 | 22.41 | 26.55 |
| + SFT | 26.00 | 26.64 | 25.17 | 20.74 | 24.64 |
| + Vanilla GRPO | 45.24 | 42.12 | 32.93 | 31.52 | 37.95 |
| + DrGRPO | 42.03 | 43.45 | 36.79 | 38.69 | 40.24 |
| + Perplexity-aware GRPO | 48.05 | 43.91 | 40.31 | 36.68 | 42.24 |
| LLaMA3.1-8B-Instruct | 24.68 | 12.88 | 11.10 | 12.45 | 15.28 |
| + SFT | 31.00 | 29.12 | 27.27 | 23.46 | 27.71 |
| + Vanilla GRPO | 57.88 | 43.38 | 25.33 | 29.06 | 38.91 |
| + DrGRPO | 55.60 | 40.64 | 31.80 | 33.67 | 40.43 |
| + Perplexity-aware GRPO | 59.82 | 41.32 | 33.96 | 35.15 | 42.56 |
| Mistral-7B-Instruct-v0.3 | 15.65 | 18.18 | 10.34 | 10.64 | 13.70 |
| + SFT | 16.85 | 23.47 | 22.27 | 22.16 | 21.19 |
| + Vanilla GRPO | 13.40 | 34.77 | 31.33 | 28.85 | 27.09 |
| + DrGRPO | 12.73 | 37.64 | 32.08 | 29.88 | 28.08 |
| + Perplexity-aware GRPO | 19.09 | 35.57 | 28.48 | 30.66 | 28.45 |

Table 3: Evaluation results on ProcessBench. GM, MT, OB, and OM correspond to the abbreviations of GSM8K, MATH, OlympiadBench, and Omni-MATH, respectively.

**Comparison of Exploration Behavior**

To illustrate exploration behavior changes caused by perplexity modulation during RL training, we present in Figure 3 the proportion curves of 'correct' and 'wrong' predictions across low, median, and high perplexity groups. Consistent with previous analyses, vanilla GRPO explores fewer trajectories that identify high-perplexity solutions as correct (as indicated by the green line with the lowest proportion of 'correct' predictions) or low-perplexity solutions as wrong (blue line with the lowest proportion of 'wrong' predictions). In contrast, our perplexity-aware GRPO clearly

| Model | LLaMA3.1-8B-Instruct-Subset | | | | Mistral-7B-Instruct-v0.3-Subset | | | | Qwen2-7B-Instruct-Subset | | | | Average | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | FPR | FNR | BI | Acc | FPR | FNR | BI | Acc | FPR | FNR | BI | Acc | $\lvert BI \rvert \downarrow$ | $Acc \uparrow$ |
| Perplexity-aware GRPO | 14.29 | 27.62 | -13.33 | 79.05 | 5.71 | 33.97 | -28.25 | 80.16 | 14.92 | 24.76 | -9.84 | 80.16 | _17.14_ | **79.79** |
| -w/o class-level aggregation | 7.62 | 27.94 | -20.32 | 82.22 | 4.76 | 42.22 | -37.46 | 76.51 | 15.87 | 30.48 | -14.60 | 76.83 | _24.13_ | _78.52_ |
| -w/o perplexity modulation | 24.13 | 20.00 | 4.13 | 77.94 | 14.29 | 34.29 | -20.00 | 75.71 | 31.75 | 20.95 | 10.79 | 73.65 | **11.64** | 75.77 |
| -w/o both | 7.94 | 39.05 | -31.11 | 76.51 | 3.17 | 54.92 | -51.75 | 70.95 | 11.75 | 34.60 | -22.86 | 76.83 | _35.24_ | _74.76_ |

Table 4: Results of the ablation study on the OPS benchmark, with Qwen2-7B-Instruct serving as the base model.
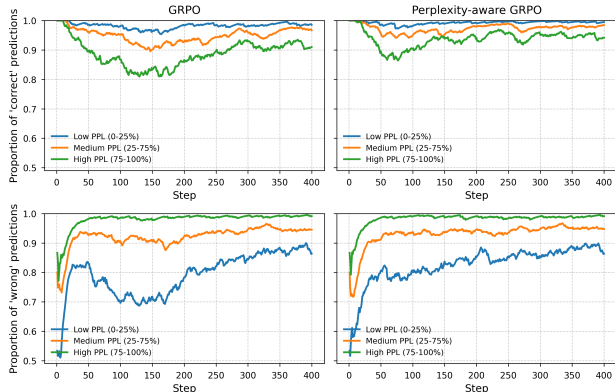


Figure 3: Exploration comparison between vanilla and perplexity-aware GRPO. Closer curves across perplexity bins imply more balanced exploration.

mitigates this preference, resulting in more balanced exploration behavior, as evidenced by the closer curves for both 'correct' and 'wrong' predictions across perplexity bins.

## Related Work

### Critiquing Capability of LLMs

The critiquing capability of LLMs is crucial for enabling scalable annotation of supervised mathematical reasoning demonstrations (Wang et al. 2024a; Lu et al. 2024a; Zhang et al. 2025b,a; Yu et al. 2025b; Gao et al. 2025), which in turn substantially enhances the multi-step mathematical reasoning performance of LLMs. To assess the critique abilities of current models, several critic-oriented benchmarks have been introduced (Lin et al. 2024; Zheng et al. 2025), consistently revealing that existing LLMs still exhibit unsatisfactory critique performance. To obtain annotated critique data for improving LLMs' critiquing skills, monte carlo sampling has been employed to efficiently annotate mathematical reasoning processes, estimating correctness based on the success rate of subsequent completions (Wang et al. 2024b; Zhang et al. 2025b). Despite these advances, they still focus on binary correctness judgments for reasoning processes, without providing more detailed or nuanced critiques. To address this limitation and enhance the interpretability of judgments, recent research has increasingly focused on annotating the critiquing process itself (Gao et al. 2025; Zhang et al. 2025a; Yu et al. 2025b). For example, Gao et al. (2025) employ GPT-4o to generate step-by-step critique annotations, providing both correctness assessments for each reasoning step and detailed explanations. Collectively, these studies aim to create high-quality supervised fine-tuning datasets to train models with improved critiquing capability. However, they pay little attention to delving into the underlying reason for the poor critiquing performance of LLMs.

### Enhancing LLM Reasoning via RL

Recent studies have demonstrated that RL is a promising approach for unlocking the long-chain, multi-step reasoning capabilities of LLMs. Such advanced reasoning skills are also essential when LLMs are tasked with performing critiques. To reduce training costs, Shao et al. (2024b) propose Group Relative Policy Optimization (GRPO), which replaces the critic model with a baseline estimated from group scores, thereby enabling large-scale RL. However, Liu et al. (2025b) identify optimization biases inherent in GRPO, such as response-level length bias and question-level difficulty bias. To address these issues, they propose Dr.GRPO, which eliminates response-length normalization and standard deviation normalization, thereby achieving a more unbiased optimization process. Furthermore, several studies (Cui et al. 2025; Yu et al. 2025a) highlight the critical role of entropy in RL, which encourages the policy to explore more diverse trajectories. In parallel, Xie et al. (2025) investigate LLMs' critiquing capability in code generation tasks and demonstrate that training critic models via RL notably enhances the critiquing capability in this domain. These advances show new potential for enhancing the ability of LLMs to critique mathematical reasoning.

## Conclusion

In this paper, we delve into and quantify the underlying reason (i.e., imbalanced evaluation preference) for the poor critiquing performance of LLMs by meticulously constructing a one-to-many problem solution benchmark. To investigate the behavior difference in depth, we conduct a statistical preference analysis oriented on perplexity and find an intriguing phenomenon—LLMs incline to judge solutions with lower perplexity as correct, which is dubbed as the imbalanced evaluation preference. Motivated by the analysis of the potential reason, we propose a novel perplexity-aware reinforcement learning algorithm to rectify the evaluation preference, which supports the LLMs to explore counter-preference trajectories that judge lower perplexity as wrong and higher perplexity as correct. Extensive experimental results on our built OPS and existing available critic benchmarks demonstrate the validity of our method.

## Acknowledgments

## References

Cui, G.; Zhang, Y.; Chen, J.; Yuan, L.; Wang, Z.; Zuo, Y.; Li, H.; Fan, Y.; Chen, H.; Chen, W.; Liu, Z.; Peng, H.; Bai, L.; Ouyang, W.; Cheng, Y.; Zhou, B.; and Ding, N. 2025. The Entropy Mechanism of Reinforcement Learning for Reasoning Language Models. *CoRR*, abs/2505.22617.

Gao, B.; Cai, Z.; Xu, R.; Wang, P.; Zheng, C.; Lin, R.; Lu, K.; Liu, D.; Zhou, C.; Xiao, W.; Liu, T.; and Chang, B. 2025. LLM Critics Help Catch Bugs in Mathematics: Towards a Better Mathematical Verifier with Natural Language Feedback. In Che, W.; Nabende, J.; Shutova, E.; and Pilehvar, M. T., eds., *Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, 14588–14604. Association for Computational Linguistics.

Grattafiori, A.; Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; Al-Dahle, A.; Letman, A.; Mathur, A.; Schelten, A.; Vaughan, A.; et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Hendrycks, D.; Burns, C.; Kadavath, S.; Arora, A.; Basart, S.; Tang, E.; Song, D.; and Steinhardt, J. 2021. Measuring Mathematical Problem Solving With the MATH Dataset. In Vanschoren, J.; and Yeung, S., eds., *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*.

Jia, W.; Jin, L.; Wei, K.; Shang, Y.; Liu, N.; Lu, Z.; Liu, Q.; Zhang, L.; Zhong, J.; and Hu, Y. 2025. U-MERE: Unconstrained Multimodal Entity and Relation Extraction with Collaborative Modeling and Order-Sensitive Optimization. In *Proceedings of the 33rd ACM International Conference on Multimedia*, MM '25, 4349–4358. New York, NY, USA: Association for Computing Machinery. ISBN 9798400720352.

Jiang, A. Q.; Sablayrolles, A.; Mensch, A.; Bamford, C.; Chaplot, D. S.; de las Casas, D.; Bressand, F.; Lengyel, G.; Lample, G.; Saulnier, L.; Lavaud, L. R.; Lachaux, M.-A.; Stock, P.; Scao, T. L.; Lavril, T.; Wang, T.; Lacroix, T.; and Sayed, W. E. 2023. Mistral 7B. arXiv:2310.06825.

Kwon, W.; Li, Z.; Zhuang, S.; Sheng, Y.; Zheng, L.; Yu, C. H.; Gonzalez, J. E.; Zhang, H.; and Stoica, I. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. arXiv:2309.06180.

Lin, Z.; Gou, Z.; Liang, T.; Luo, R.; Liu, H.; and Yang, Y. 2024. CriticBench: Benchmarking LLMs for Critique-Correct Reasoning. In Ku, L.; Martins, A.; and Srikumar, V., eds., *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, 1552–1587. Association for Computational Linguistics.

Liu, N.; Zhu, J.; Ma, Y.; Lu, Z.; Xu, W.; Yang, Y.; Zhong, J.; and Wei, K. 2025a. SARA: Salience-Aware Reinforced Adaptive Decoding for Large Language Models in Abstractive Summarization. In Che, W.; Nabende, J.; Shutova, E.; and Pilehvar, M. T., eds., *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 25450–25463. Vienna, Austria: Association for Computational Linguistics. ISBN 979-8-89176-251-0.

Liu, Z.; Chen, C.; Li, W.; Qi, P.; Pang, T.; Du, C.; Lee, W. S.; and Lin, M. 2025b. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*.

Lu, Z.; Jin, L.; Chen, Z.; Tian, C.; Sun, X.; Li, X.; Zhang, Y.; Li, Q.; and Xu, G. 2024a. Relation-Aware Multi-Pass Comparison Deconfounded Network for Change Captioning. *IEEE Transactions on Circuits and Systems for Video Technology*, 34(12): 13349–13363.

Lu, Z.; Jin, L.; Li, P.; Tian, Y.; Zhang, L.; Wang, S.; Xu, G.; Tian, C.; and Cai, X. 2024b. Rethinking the Reversal Curse of LLMs: a Prescription from Human Knowledge Reversal. In Al-Onaizan, Y.; Bansal, M.; and Chen, Y.-N., eds., *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 7518–7530. Miami, Florida, USA: Association for Computational Linguistics.

Lu, Z.; Jin, L.; Xu, G.; Hu, L.; Liu, N.; Li, X.; Sun, X.; Zhang, Z.; and Wei, K. 2023. Narrative Order Aware Story Generation via Bidirectional Pretraining Model with Optimal Transport Reward. In Bouamor, H.; Pino, J.; and Bali, K., eds., *Findings of the Association for Computational Linguistics: EMNLP 2023*, 6274–6287. Singapore: Association for Computational Linguistics.

Lu, Z.; Tian, C.; PeiguangLi, P.; Jin, L.; Wang, S.; Jia, W.; Shen, Y.; and Xu, G. 2025. PIPER: Benchmarking and Prompting Event Reasoning Boundary of LLMs via Debiasing-Distillation Enhanced Tuning. In Che, W.; Nabende, J.; Shutova, E.; and Pilehvar, M. T., eds., *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 28591–28613. Vienna, Austria: Association for Computational Linguistics. ISBN 979-8-89176-251-0.

Qwen; :; Yang, A.; Yang, B.; Zhang, B.; Hui, B.; Zheng, B.; Yu, B.; Li, C.; Liu, D.; Huang, F.; Wei, H.; Lin, H.; Yang, J.; Tu, J.; Zhang, J.; Yang, J.; Yang, J.; Zhou, J.; Lin, J.; Dang, K.; Lu, K.; Bao, K.; Yang, K.; Yu, L.; Li, M.; Xue, M.; Zhang, P.; Zhu, Q.; Men, R.; Lin, R.; Li, T.; Tang, T.; Xia, T.; Ren, X.; Ren, X.; Fan, Y.; Su, Y.; Zhang, Y.; Wan, Y.; Liu, Y.; Cui, Z.; Zhang, Z.; and Qiu, Z. 2025. Qwen2.5 Technical Report. arXiv:2412.15115.

Shao, Z.; Wang, P.; Zhu, Q.; Xu, R.; Song, J.; Zhang, M.; Li, Y. K.; Wu, Y.; and Guo, D. 2024. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. *CoRR*, abs/2402.03300.

Sheng, G.; Zhang, C.; Ye, Z.; Wu, X.; Zhang, W.; Zhang, R.; Peng, Y.; Lin, H.; and Wu, C. 2025. HybridFlow: A Flexible and Efficient RLHF Framework. In *Proceedings of the Twentieth European Conference on Computer Systems, EuroSys 2025, Rotterdam, The Netherlands, 30 March 2025 - 3 April 2025*, 1279–1297. ACM.

Team, Q. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

Tian, C.; Lu, Z.; Zhang, Z.; Yang, H.; Cao, W.; Guo, Z.; Sun, X.; and Jin, L. 2025. HyperMixer: Specializable Hypergraph Channel Mixing for Long-term Multivariate Time Series Forecasting. In Walsh, T.; Shah, J.; and Kolter, Z., eds., *AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025, Philadelphia, PA, USA*, 20885–20893. AAAI Press.

Wang, P.; Li, L.; Shao, Z.; Xu, R.; Dai, D.; Li, Y.; Chen, D.; Wu, Y.; and Sui, Z. 2024a. Math-Shepherd: Verify and Reinforce LLMs Step-by-step without Human Annotations. In Ku, L.-W.; Martins, A.; and Srikumar, V., eds., *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 9426–9439. Bangkok, Thailand: Association for Computational Linguistics.

Wang, P.; Li, L.; Shao, Z.; Xu, R.; Dai, D.; Li, Y.; Chen, D.; Wu, Y.; and Sui, Z. 2024b. Math-Shepherd: Verify and Reinforce LLMs Step-by-step without Human Annotations. In Ku, L.; Martins, A.; and Srikumar, V., eds., *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, 9426–9439. Association for Computational Linguistics.

Xie, Z.; Chen, L.; Mao, W.; Xu, J.; Kong, L.; et al. 2025. Teaching Language Models to Critique via Reinforcement Learning. *arXiv preprint arXiv:2502.03492*.

Yang, A.; Li, A.; Yang, B.; Zhang, B.; Hui, B.; Zheng, B.; Yu, B.; Gao, C.; Huang, C.; Lv, C.; et al. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.

Yu, Q.; Zhang, Z.; Zhu, R.; Yuan, Y.; Zuo, X.; Yue, Y.; Fan, T.; Liu, G.; Liu, L.; Liu, X.; Lin, H.; Lin, Z.; Ma, B.; Sheng, G.; Tong, Y.; Zhang, C.; Zhang, M.; Zhang, W.; Zhu, H.; Zhu, J.; Chen, J.; Chen, J.; Wang, C.; Yu, H.; Dai, W.; Song, Y.; Wei, X.; Zhou, H.; Liu, J.; Ma, W.; Zhang, Y.; Yan, L.; Qiao, M.; Wu, Y.; and Wang, M. 2025a. DAPO: An Open-Source LLM Reinforcement Learning System at Scale. *CoRR*, abs/2503.14476.

Yu, Y.; Chen, Z.; Zhang, A.; Tan, L.; Zhu, C.; Pang, R. Y.; Qian, Y.; Wang, X.; Gururangan, S.; Zhang, C.; Kambadur, M.; Mahajan, D.; and Hou, R. 2025b. Self-Generated Critiques Boost Reward Modeling for Language Models. In Chiruzzo, L.; Ritter, A.; and Wang, L., eds., *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)*, 11499–11514. Albuquerque, New Mexico: Association for Computational Linguistics. ISBN 979-8-89176-189-6.

Zhang, L.; Hosseini, A.; Bansal, H.; Kazemi, M.; Kumar, A.; and Agarwal, R. 2025a. Generative Verifiers: Reward Modeling as Next-Token Prediction. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.

Zhang, Z.; Zheng, C.; Wu, Y.; Zhang, B.; Lin, R.; Yu, B.; Liu, D.; Zhou, J.; and Lin, J. 2025b. The Lessons of Developing Process Reward Models in Mathematical Reasoning.

In Che, W.; Nabende, J.; Shutova, E.; and Pilehvar, M. T., eds., *Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, 10495–10516. Association for Computational Linguistics.

Zheng, C.; Zhang, Z.; Zhang, B.; Lin, R.; Lu, K.; Yu, B.; Liu, D.; Zhou, J.; and Lin, J. 2025. ProcessBench: Identifying Process Errors in Mathematical Reasoning. arXiv:2412.06559.