



# Harder Is Better: Boosting Mathematical Reasoning via Difficulty-Aware GRPO and Multi-Aspect Question Reformulation

Yanqi Dai (代彦琪)<sup>12</sup>, Yuxiang Ji<sup>23</sup>, Xiao Zhang<sup>24</sup>, Yong Wang<sup>2\*</sup>, Xiangxiang Chu<sup>2</sup>, Zhiwu Lu<sup>1\*</sup>

<sup>1</sup>Gaoling School of Artificial Intelligence, Renmin University of China

<sup>2</sup>AMAP, Alibaba Group   <sup>3</sup>Xiamen University   <sup>4</sup>Dalian University of Technology

**Motivation:** Reinforcement Learning with Verifiable Rewards (RLVR)

- **[Algorithm] GRPO:** implicitly weaken the update magnitude of harder questions.
- **[Data] WizardMath, MetaMath:** need to generate both questions and answers from scratch, or just rephrase questions to enhance diversity rather than difficulty.

**Preliminaries:** Optimization Objective of Group Relative Policy Optimization (GRPO)

Specifically, GRPO optimizes the policy model  $\pi_\theta$  by maximizing the following objective:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E} \left[ q \sim \mathcal{D}, \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | q) \right] \frac{1}{\sum_{i=1}^G |o_i|} \sum_{i=1}^G \sum_{t=1}^{|o_i|} \left\{ \min \left[ I_{it}(\theta) \hat{A}_{\text{GR},i}, \text{clip} \left( I_{it}(\theta), 1 - \varepsilon, 1 + \varepsilon \right) \hat{A}_{\text{GR},i} \right] \right\}, \quad (1)$$

$$\text{where } I_{it}(\theta) = \frac{\pi_\theta(o_{i,t} | q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t} | q, o_{i,<t})}, \hat{A}_{\text{GR},i} = \frac{r_i - \text{mean}(\{r_i\}_{i=1}^G)}{\text{std}(\{r_i\}_{i=1}^G)}. \quad (2)$$

## Optimization Objective of Difficulty-Aware Group Policy Optimization (DGPO)

Specifically, the optimization objective of DGPO is defined as follows:

$$\mathcal{J}_{\text{DGPO}}(\theta) = \mathbb{E} \left[ \{q_s\}_{s=1}^B \sim \mathcal{D}, \{o_{si}\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | q_s) \right]$$

$$\frac{1}{\sum_{s=1}^{B_v} \sum_{i=1}^G |o_{si}|} \sum_{s=1}^{B_v} \lambda_s \sum_{i=1}^G \sum_{t=1}^{|o_{si}|} \left\{ \min \left[ I_{sit}(\theta) \hat{A}_{\text{DG},si}, \text{clip} \left( I_{sit}(\theta), 1 - \varepsilon, 1 + \varepsilon \right) \hat{A}_{\text{DG},si} \right] \right\}, \quad (3)$$

where  $I_{sit}(\theta)$  is the importance sampling ratio of the token  $o_{si,t}$ , and  $\hat{A}_{\text{DG},si}$  is the advantage of the response  $o_i$  obtained by DGAE, respectively given by:

$$I_{sit}(\theta) = \frac{\pi_{\theta}(o_{si,t} | q_s, o_{si,<t})}{\pi_{\theta_{\text{old}}}(o_{si,t} | q_s, o_{si,<t})}, \quad \hat{A}_{\text{DG},si} = \frac{r_{si} - \text{mean}(\{r_{si}\}_{i=1}^G)}{\text{MAD}(\{r_{si}\}_{i=1}^G)}, \quad (4)$$

$$\text{where } \text{MAD}(\{r_{si}\}_{i=1}^G) = \frac{1}{G} \sum_{i=1}^G |r_{si} - \text{mean}(\{r_{si}\}_{i=1}^G)|. \quad (5)$$

Here,  $\text{MAD}(\cdot)$  denotes the mean absolute deviation function. Furthermore,  $\lambda_s$  is the difficulty-aware weight for the query  $q_s$  computed by DQW as follows:

$$\lambda_s = B_v \cdot \frac{\exp(D_s/T)}{\sum_{s=1}^{B_v} \exp(D_s/T)}, \quad \text{where } D_s = -\text{mean}(\{r_{si}\}_{i=1}^G). \quad (6)$$

## Optimization Objective of Difficulty-Aware Group Policy Optimization (DGPO) Difficulty-balanced Group Advantage Estimation (DGAE)

Consider a single question  $q$  and its corresponding responses  $\{o_i\}_{i=1}^G$ , the unclipped policy gradient calculated in GRPO is as follows:

$$\begin{aligned}
 g_{\text{GRPO}} &= \frac{1}{\sum_{i=1}^G |o_i|} \sum_{i=1}^G \sum_{t=1}^{|o_i|} \hat{A}_{\text{GR},i} \nabla_{\theta} I_{it}(\theta) \\
 &= \frac{1}{\sum_{i=1}^G |o_i|} \sum_{i=1}^G \sum_{t=1}^{|o_i|} \underbrace{\text{sgn}(\hat{A}_{\text{GR},i})}_{\text{update direction}} \underbrace{|\hat{A}_{\text{GR},i}|}_{\text{update magnitude}} \underbrace{\text{detach}(I_{it}(\theta))}_{\text{importance sampling ratio}} \underbrace{\nabla_{\theta} \log(\pi_{\theta}(o_{i,t} | q, o_{i,<t}))}_{\text{likelihood gradient}}, \quad (7)
 \end{aligned}$$

**Theorem 1** (Update Magnitude for a Single Question using GRAE). *Given a single question  $q$  and its corresponding responses  $\{o_i\}_{i=1}^G$ , each query-response pair receives a binary accuracy reward  $r_i \in \{0, 1\}$ , and  $p$  represents the accuracy rate, i.e., the proportion for a reward of 1. Then, the total update magnitude without clipping for the single question  $q$  when using GRAE satisfies:*

$$\sum_{i=1}^G |\hat{A}_{\text{GR},i}| = \sum_{i=1}^G \left| \frac{r_i - \text{mean}(\{r_i\}_{i=1}^G)}{\text{std}(\{r_i\}_{i=1}^G)} \right| = 2G\sqrt{p(1-p)}, \text{ where } p = \frac{1}{G} \sum_{i=1}^G r_i. \quad (8)$$

*This total update magnitude varies with respect to the accuracy rate  $p$ , reaching its maximum when  $p = 0.5$  and gradually decreasing as  $p$  approaches either 0 or 1.*

## Optimization Objective of Difficulty-Aware Group Policy Optimization (DGPO)

### Difficulty-balanced Group Advantage Estimation (DGAE)

Specifically, the advantage function of DGAE is defined as follows:

$$\hat{A}_{\text{DG},i} = \frac{r_i - \text{mean}(\{r_i\}_{i=1}^G)}{\text{MAD}(\{r_i\}_{i=1}^G)}, \quad \text{where } \text{MAD}(\{r_i\}_{i=1}^G) = \frac{1}{G} \sum_{i=1}^G |r_i - \text{mean}(\{r_i\}_{i=1}^G)|. \quad (9)$$

**Theorem 2** (Update Magnitude for a Single Question using DGAE). *Given a single question  $q$  and its corresponding responses  $\{o_i\}_{i=1}^G$ , each query-response pair receives a reward  $r_i$ . Then, the total update magnitude without clipping for the single question  $q$  when using DGAE satisfies:*

$$\sum_{i=1}^G |\hat{A}_{\text{DG},i}| = \sum_{i=1}^G \left| \frac{r_i - \text{mean}(\{r_i\}_{i=1}^G)}{\frac{1}{G} \sum_{i=1}^G |r_i - \text{mean}(\{r_i\}_{i=1}^G)|} \right| = G. \quad (10)$$

## Optimization Objective of Difficulty-Aware Group Policy Optimization (DGPO)

### Difficulty-aware Question-level Weighting (DQW)

Core idea: **prioritizes learning from more challenging yet solvable questions**

Specifically, DQW assigns a weight  $\lambda_s$  to each question  $q_s$  as follows:

$$\lambda_s = B_v \cdot \frac{\exp(D_s/T)}{\sum_{s=1}^{B_v} \exp(D_s/T)}, \text{ where } D_s = -\text{mean}(\{r_{si}\}_{i=1}^G). \quad (11)$$

Here,  $D_s$  is the negative mean reward across all responses of the question  $q_s$ , serving as a measure of its relative difficulty at the current training stage.

**DGPO = DGAE + DQW**  
**“balance-then-reweight”**

## Multi-Aspect Question Reformulation (MQR)

### Core Instructions for Multi-Aspect Question Reformulation

1. **Background:** Add a story background that is not related to the core mathematical content of the given question, but seems to be related to the question. If the given question already has such a background, change it to a new, complexer background.
2. **Term:** Invent a new, abstract mathematical term to define a concept that is central to the given question, and restate the entire question using this term.
3. **Sub-Problem:** Convert a key numerical condition of the given question which have a definite value into an independent sub-problem. The sub-problem may belong to any branch of mathematics (e.g., algebra, geometry, number theory, combinatorics).

A critical constraint is that *all reformulations must preserve the original gold answer*.

The newly generated questions respectively challenge the policy model's ability to:

1. identify critical mathematical information amidst noise;
2. grasp abstract mathematical concepts; and
3. perform reasoning that requires multiple steps and cross-domain knowledge.

## Main Results

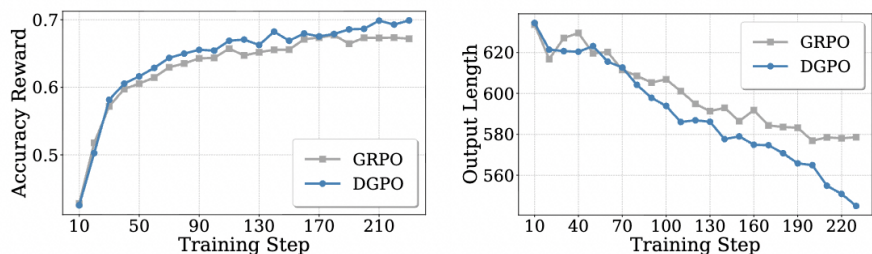
Table 1: Comparative results of methods trained on the MATH dataset using Qwen2.5-Math-7B.

Methods	AIME24	AIME25	AMC23	MATH500	Minerva	Olympiad	Avg./ $\Delta_{GRPO}$
Base Model	12.19	4.79	35.23	48.60	15.07	16.33	22.04
GRPO	20.94	8.44	58.98	72.20	27.76	37.33	37.61
Dr.GRPO	21.04	8.23	58.59	72.05	28.58	35.89	37.40 <sup>-0.21</sup>
GPG	21.98	9.06	59.61	72.05	27.21	37.67	37.93 <sup>+0.32</sup>
DAPO	21.25	8.75	58.20	72.70	29.50	37.22	37.94 <sup>+0.33</sup>
GSPO	19.38	8.33	<u>60.16</u>	73.00	28.12	37.26	37.71 <sup>+0.10</sup>
GRPO-AD	21.56	9.48	59.06	73.25	29.14	37.07	38.26 <sup>+0.65</sup>
DGPO	23.85	10.21	<b>61.02</b>	74.25	31.07	38.33	39.79 <sup>+2.18</sup>
MQR	<b>25.00</b>	<u>11.77</u>	59.38	<u>77.85</u>	<u>31.43</u>	<u>40.81</u>	<u>41.04</u> <sup>+3.43</sup>
MathForge	<u>24.58</u>	<b>12.60</b>	59.84	<b>79.95</b>	<b>33.36</b>	<b>42.67</b>	<b>42.17</b> <sup>+4.56</sup>

Table 2: Comparative results of methods trained on the MATH dataset using varying base models.

Methods	AIME24	AIME25	AMC23	MATH500	Minerva	Olympiad	Avg./ $\Delta_{GRPO}$
Qwen2.5-Math-1.5B	6.87	3.65	30.94	34.95	8.55	21.93	17.82
+ GRPO	11.35	3.96	46.48	64.85	20.13	29.59	29.39
+ DGPO	11.25	<u>5.73</u>	49.84	65.45	21.14	30.85	30.71 <sup>+1.32</sup>
+ MQR	<u>11.98</u>	5.42	<u>50.08</u>	<u>69.65</u>	<u>23.81</u>	<u>33.67</u>	<u>32.44</u> <sup>+3.05</sup>
+ MathForge	<b>13.23</b>	<b>7.71</b>	<b>52.34</b>	<b>70.10</b>	<b>25.74</b>	<b>33.89</b>	<b>33.84</b> <sup>+4.45</sup>
Qwen2.5-3B	2.81	0.73	22.66	48.65	13.69	19.37	17.99
+ GRPO	5.31	<u>1.56</u>	33.28	63.35	22.89	26.41	25.47
+ DGPO	<b>6.98</b>	<u>1.56</u>	36.56	<b>65.80</b>	25.28	26.96	27.19 <sup>+1.72</sup>
+ MQR	5.10	<u>1.56</u>	<u>39.53</u>	65.20	<u>25.74</u>	<u>29.19</u>	<u>27.72</u> <sup>+2.25</sup>
+ MathForge	<u>5.73</u>	<b>1.77</b>	<b>40.70</b>	<u>65.40</u>	<b>28.86</b>	<b>31.59</b>	<b>29.01</b> <sup>+3.54</sup>
DeepSeek-Math-7B	0.42	0.10	13.28	31.05	9.56	9.00	10.57
+ GRPO	0.63	0.10	19.14	41.45	14.71	13.44	14.91
+ DGPO	<u>1.98</u>	0.42	<u>21.02</u>	41.85	<u>18.93</u>	15.00	16.53 <sup>+1.62</sup>
+ MQR	<u>1.98</u>	<b>0.83</b>	20.86	<b>44.25</b>	17.00	<u>15.74</u>	<u>16.78</u> <sup>+1.87</sup>
+ MathForge	<b>3.12</b>	<u>0.73</u>	<b>21.72</b>	<u>43.60</u>	<b>20.68</b>	<b>16.74</b>	<b>17.77</b> <sup>+2.86</sup>

## Analysis of DGPO



(a) Accuracy Reward.

(b) Output Length.

Figure 1: Training dynamics of DGPO vs. GRPO evaluated on the MATH500 benchmark. Both models are trained on MATH using Qwen2.5-Math-7B.

Table 3: Ablation Results of DGPO trained on the MATH dataset using Qwen2.5-Math-7B.

Methods	AIME24	AIME25	AMC23	MATH500	Minerva	Olympiad	Avg./ $\Delta_{GRPO}$
GRPO	20.94	8.44	58.98	72.20	27.76	37.33	37.61
DGPO (w/o DGAE & DQW)	20.21	9.06	59.45	72.40	28.58	36.56	37.71 <sup>+0.10</sup>
DGPO (w/o DQW)	<u>21.77</u>	<u>9.69</u>	<u>60.00</u>	<u>73.45</u>	<u>29.04</u>	<u>37.93</u>	<u>38.65</u> <sup>+1.04</sup>
DGPO (full)	<b>23.85</b>	<b>10.21</b>	<b>61.02</b>	<b>74.25</b>	<b>31.07</b>	<b>38.33</b>	<b>39.79</b> <sup>+2.18</sup>
DGPO ( $T = 1.0$ )	<u>23.12</u>	9.06	59.45	74.15	<u>30.61</u>	37.78	39.03 <sup>+1.42</sup>
DGPO ( $T = 2.0$ )	<b>23.85</b>	<u>10.21</u>	<u>61.02</u>	<u>74.25</u>	<b>31.07</b>	<b>38.33</b>	<b>39.79</b> <sup>+2.18</sup>
DGPO ( $T = 5.0$ )	22.81	<b>11.35</b>	60.55	73.80	30.42	<u>38.26</u>	<u>39.53</u> <sup>+1.92</sup>
DGPO ( $T = 10.0$ )	21.35	9.79	<b>62.27</b>	<b>74.55</b>	29.96	37.67	39.27 <sup>+1.66</sup>

Table 4: Synergistic results of DGPO with other policy optimization methods trained on the MATH dataset using Qwen2.5-Math-7B.

Methods	AIME24	AIME25	AMC23	MATH500	Minerva	Olympiad	Average
GPG	<b>21.98</b>	9.06	59.61	72.05	27.21	37.67	37.93
+ DGPO	21.77	<b>10.00</b>	<b>60.00</b>	<b>73.45</b>	<b>30.06</b>	<b>38.26</b>	<b>38.92</b>
DAPO	21.25	8.75	58.20	72.70	29.50	37.22	37.94
+ DGPO	<b>24.48</b>	<b>9.79</b>	<b>58.75</b>	<b>74.90</b>	<b>31.99</b>	<b>39.56</b>	<b>39.91</b>
GSPO	19.38	8.33	<b>60.16</b>	73.00	28.12	37.26	37.71
+ DGPO	<b>23.33</b>	<b>10.00</b>	59.14	<b>74.15</b>	<b>30.88</b>	<b>38.41</b>	<b>39.32</b>

Table 5: Comparative results of methods trained on the GEOQA-8k dataset using Qwen2.5-VL-3B-Instruct in the multimodal domain.

Methods	Base Model	GRPO	Dr.GRPO	GPG	DAPO	GSPO	GRPO-AD	DGPO
GeoQA/ $\Delta_{GRPO}$	39.79	57.43	57.96 <sup>+0.53</sup>	<u>59.02</u> <sup>+1.59</sup>	<u>59.02</u> <sup>+1.59</sup>	57.16 <sup>-0.27</sup>	58.09 <sup>+0.66</sup>	<b>59.95</b> <sup>+2.52</sup>

## Analysis of DQW

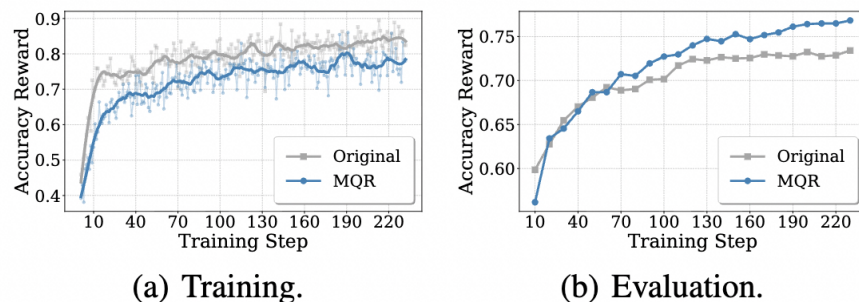


Figure 2: Training dynamics of Original vs. MQR on training and evaluation data. Both models are trained on MATH and evaluated on MATH500 using Qwen2.5-Math-7B.

Table 6: Comparative results of methods trained on the original data vs. the MQR-augmented data using DGPO and varying base models.

Models	Data	AIME24	AIME25	AMC23	MATH500	Minerva	Olympiad	Average
Qwen2.5-Math-7B	Ori.	<b>26.46</b>	9.17	58.67	74.65	31.62	38.81	39.90
	MQR	24.58	<b>12.60</b>	<b>59.84</b>	<b>79.95</b>	<b>33.36</b>	<b>42.67</b>	<b>42.17</b>
Qwen2.5-Math-1.5B	Ori.	11.98	5.21	50.62	68.40	24.26	32.59	32.18
	MQR	<b>13.23</b>	<b>7.71</b>	<b>52.34</b>	<b>70.10</b>	<b>25.74</b>	<b>33.89</b>	<b>33.84</b>
Qwen2.5-3B	Ori.	<b>6.04</b>	1.35	37.66	65.05	25.28	27.93	27.22
	MQR	5.73	<b>1.77</b>	<b>40.70</b>	<b>65.40</b>	<b>28.86</b>	<b>31.59</b>	<b>29.01</b>
DeepSeek-Math-7B	Ori.	2.19	0.21	21.02	<b>43.60</b>	18.29	14.52	16.64
	MQR	<b>3.12</b>	<b>0.73</b>	<b>21.72</b>	<b>43.60</b>	<b>20.68</b>	<b>16.74</b>	<b>17.77</b>

Table 7: Ablation Results of MQR on the MATH dataset using Qwen2.5-Math-7B.

Data	AIME24	AIME25	AMC23	MATH500	Minerva	Olympiad	Avg./ $\Delta_{Ori.}$
Original	<u>26.46</u>	9.17	58.67	74.65	31.62	38.81	39.90
MetaMath-Rephrasing	25.21	<u>11.35</u>	<u>59.45</u>	76.70	31.71	39.93	40.73 <sup>+0.83</sup>
Original + Background	25.52	10.73	58.59	77.50	32.90	40.48	40.95 <sup>+1.05</sup>
Original + Term	25.52	11.15	58.98	<u>77.75</u>	33.09	40.93	41.24 <sup>+1.34</sup>
Original + Sub-Problem	<b>26.67</b>	10.94	58.75	77.05	<b>34.38</b>	<u>41.36</u>	<u>41.53</u> <sup>+1.63</sup>
MQR	24.58	<b>12.60</b>	<b>59.84</b>	<b>79.95</b>	<u>33.36</u>	<b>42.67</b>	<b>42.17</b> <sup>+2.27</sup>

Table 8: Comparative results of MQR using varying reformulator models on the MATH dataset.

Reformulators	AIME24	AIME25	AMC23	MATH500	Minerva	Olympiad	Avg./ $\Delta_{Ori.}$
Original	26.46	9.17	58.67	74.65	31.62	38.81	39.90
Qwen2.5-7B-Instruct	<u>25.10</u>	11.98	58.67	76.85	33.00	40.96	41.09 <sup>+1.19</sup>
Qwen3-30B-A3B-Thinking	<b>25.73</b>	<u>12.29</u>	<b>59.84</b>	<u>78.85</u>	<u>33.18</u>	<u>41.22</u>	<u>41.85</u> <sup>+1.95</sup>
OpenAI o3	24.58	<b>12.60</b>	<b>59.84</b>	<b>79.95</b>	<b>33.36</b>	<b>42.67</b>	<b>42.17</b> <sup>+2.27</sup>



Yechiel  
Singapore



Scan QR code to add me



Yechiel  
WhatsApp contact



# Thank you for your attention!

---