

# Good Reasoning Makes Good Demonstrations: Implicit Reasoning Quality Supervision via In-Context Reinforcement Learning

Anonymous ACL submission

## Abstract

Reinforcement Learning with Verifiable Rewards (RLVR) improves reasoning in large language models but treats all correct solutions equally, potentially reinforcing flawed traces that get correct answers by chance. We observe that *better reasoning are better teachers*: high-quality solutions serve as more effective demonstrations than low-quality ones. We term this teaching ability **Demonstration Utility**, and show that the policy model’s own in-context learning ability provides an efficient way to measure it, yielding a quality signal termed **Evidence Gain**. To employ this signal during training, we introduce **In-Context RLVR**. By Bayesian analysis, we show that this objective implicitly reweights rewards by Evidence Gain, assigning higher weights to high-quality traces and lower weights to low-quality ones, without requiring costly computation or external evaluators. Experiments on mathematical benchmarks show improvements in both accuracy and reasoning quality over standard RLVR.

## 1 Introduction

Reinforcement Learning with Verifiable Rewards (RLVR) has emerged as a powerful paradigm for improving LLM reasoning (Shao et al., 2024; Guo et al., 2025a), especially in domains such as mathematics where correctness can be checked by rules (Su et al., 2025c). By using outcome-level supervision, RLVR avoids costly process annotations and scales well (Mroueh, 2025). However, this simplicity comes with a limitation: all correct solutions receive equal reward, regardless of the reasoning used to obtain them (Do et al., 2025). This is problematic since models can produce flawed reasoning traces that coincidentally get correct answers, particularly when final answers are simple values that can be guessed (Guo et al., 2025b). Consequently, reinforcing such traces may corrupt internal reasoning strategies, degrading performance on other problems (MacDiarmid et al., 2025).

A natural solution is to use process reward models (PRMs) (Zhang et al., 2025b; Ye et al., 2025) that score intermediate steps. However, PRMs typically require extensive human annotation or auxiliary trained evaluators (Lightman et al., 2024). This raises a key question: *Can we encourage high-quality reasoning within RLVR without requiring step-level supervision or trained reward models?*

**Demonstration Utility as Global Quality Signal.** Our key insight is that *high-quality reasoning traces are better teachers than low-quality ones* (Min et al., 2022). Consider two solutions that both arrive at the correct answer: one is coherent and complete; the other contains redundant or unclear steps. When used as in-context demonstrations, the former provides transferable problem-solving patterns that help the model generate better solutions, while the latter provides less reference value (Li et al., 2025). We term this teaching ability **Demonstration Utility**. Crucially, the policy model’s own in-context learning (ICL) ability provides a natural way to measure Demonstration Utility. Specifically, we construct a held-out validation set composed of questions and high-quality reference reasoning traces. We propose computing the average increase in the model’s log-likelihood of generating these references after a candidate reasoning trace is prepended as a demonstration. We call this measure **Evidence Gain** (§2). Unlike PRMs that require external evaluators, Evidence Gain leverages the intrinsic ICL capability of the policy model itself. Experiments in Section 2 confirm that this intrinsic signal effectively distinguishes good reasoning from bad.

**Implicit Reward Reweighting via In-Context RLVR.** While Evidence Gain provides a reasoning quality signal, computing it as rewards would introduce substantial overhead. Fortunately, we show that this explicit computation is unnecessary. Our key idea is to reverse the process: in-

083 instead of computing Evidence Gain *after* genera-  
084 tion as rewards, we use the same validation set  
085 to guide training *before* generation. Specifically,  
086 before each rollout, we sample a demonstration  
087 from the validation set and prepend it to the current  
088 question, then perform standard RL updates in this  
089 demonstration-conditioned setting, a procedure we  
090 term **In-Context RLVR**. Via Bayesian analysis  
091 (§3), we show that this training objective is equiva-  
092 lent to standard zero-shot RLVR but with *rewards*  
093 *implicitly reweighted by Evidence Gain*. Conse-  
094 quently, high-quality traces with greater teaching  
095 utility receive amplified gradient signals, while  
096 low-quality traces receive relatively lower weights  
097 through this implicit reweighting mechanism.

098 **Contributions.** (1) We introduce **Evidence Gain**,  
099 a quality signal that measures reasoning quality by  
100 leveraging the policy model’s intrinsic ICL abil-  
101 ity, requiring no external evaluators or step-level  
102 supervision. (2) We show that this signal can be  
103 seamlessly integrated into training via **In-Context**  
104 **RLVR**, which prepends demonstrations during  
105 training to upweight high Evidence Gain traces.  
106 (3) Experiments across mathematical benchmarks  
107 validate that our method improves both accuracy  
108 and reasoning quality over competitive baselines,  
109 while introducing less than 5% training overhead.

## 110 2 Evidence Gain as Quality Measure

111 This section formally defines **Evidence Gain**, a  
112 quality signal that measures reasoning quality by  
113 leveraging the policy model’s intrinsic in-context  
114 learning ability, and validates it empirically.

115 Our basic idea is that, when used as demon-  
116 strations, high-quality reasoning traces provide  
117 more valuable problem-solving patterns, while low-  
118 quality reasoning (even with correct answers) pro-  
119 vides less reference value due to flaws such as in-  
120 consistent logic (Li et al., 2025). This motivates us  
121 to quantify the quality of a solution by its teaching  
122 ability as a demonstration, formalized as follows.

123 **Definition.** Let  $\pi_\theta$  denote policy model. Given  
124 a question  $q$ , a model-generated reasoning trace  $r$ ,  
125 and a held-out validation set  $\mathcal{E} = \{(e_q, e_r)\}$  com-  
126 posed of questions  $e_q$  and high-quality reference  
127 reasoning traces  $e_r$ , we define **Evidence Gain** as:

$$128 \Delta(q, r) = \mathbb{E}_{e \sim \mathcal{E}} [\log \pi_\theta(e_r | q, r, e_q) - \log \pi_\theta(e_r | e_q)] \quad (1)$$

129 Intuitively,  $\Delta$  measures how much prepending the  
130 pair  $(q, r)$  improves the model’s ability to gener-  
131 ate reference solutions. Averaging over  $\mathcal{E}$  ensures

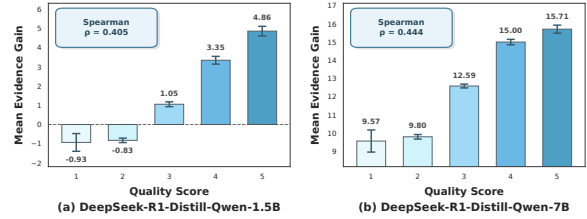


Figure 1: Mean Evidence Gain by quality score with 95% confidence intervals on two models.

132 that high  $\Delta$  reflects transferable reasoning patterns  
133 rather than spurious matches to any single sample.

134 **Empirical Validation.** We validate Evidence  
135 Gain on DeepSeek-R1-Distill-Qwen at 1.5B and  
136 7B scales using KlearReasoner-MathSub-30K  
137 dataset (Su et al., 2025a). First, we sample 3,000  
138 questions, generate 8 responses per question, and  
139 retain traces with correct final answers (12,251  
140 for 1.5B and 16,910 for 7B). We then employ  
141 DeepSeek-V3.2 (DeepSeek-AI et al., 2025), a  
142 strong LLM-based evaluator, to assess the reason-  
143 ing quality. The evaluator scores each solution  
144 across multiple dimensions including logical co-  
145 herence and redundancy, and assigns an overall  
146 quality score on a 1–5 scale. Next, we sample 100  
147 new questions to construct the validation set  $\mathcal{E}$ .  
148 For each  $e_q \in \mathcal{E}$ , we generate a correct solution  
149 using DeepSeek-R1-0528 (Guo et al., 2025a). We  
150 leverage a feature of DeepSeek-R1: it first produces  
151 a draft chain-of-thought inside `<think>` `</think>`,  
152 and then outputs a more polished reasoning so-  
153 lution afterward. We treat the content following  
154 `</think>` as the high-quality reference trace  $e_r$ .

155 Figure 1 shows an interesting patterns: model  
156 ability determines the absolute baseline of  $\Delta$ , while  
157 reasoning quality differentiates its relative magni-  
158 tude. The 7B model, with stronger ICL ability, ex-  
159 tracts useful information from any reasoning trace,  
160 resulting in uniformly positive  $\Delta$  values greatly  
161 higher than those of 1.5B. However, the relative  
162 ordering remains consistent within each scale: high-  
163 quality traces yield higher  $\Delta$  than low-quality ones.  
164 This relative difference is what matters for RL train-  
165 ing. These results confirm that Evidence Gain ef-  
166 fectively distinguishes reasoning quality. Human  
167 evaluation in Appendix C.2 supports these findings.

## 168 3 Implicit Reward Reweighting

169 The correlation between Evidence Gain ( $\Delta$ ) and  
170 reasoning quality (Section 2) suggests that up-  
171 weighting high- $\Delta$  traces during training could im-  
172 prove reasoning.

$\Delta$  for each rollout is prohibitively expensive.<sup>1</sup> In this section, we show that explicit computation is unnecessary. Evidence Gain can be seamlessly integrated into training through **In-Context RLVR**.

The core idea is to reverse the process. Instead of computing Evidence Gain *after* generation to reweight rewards, we utilize the validation set to guide training *before* generation. This follows from a Bayesian identity (derivation in Appendix E):

$$\underbrace{\pi_{\theta}(r|e, q)}_{\text{conditioned policy}} = \underbrace{\pi_{\theta}(r|q)}_{\text{base policy}} \cdot \underbrace{\frac{\pi_{\theta}(e_r|q, r, e_q)}{\pi_{\theta}(e_r|e_q)}}_{\text{likelihood ratio}},$$

where  $q$  is the query,  $r$  is the generated reasoning trace, and  $e = (e_q, e_r)$  is a demonstration sampled from  $\mathcal{E}$ . This identity shows that the conditioned policy  $\pi_{\theta}(r|e, q)$  equals the base policy  $\pi_{\theta}(r|q)$  reweighted by a likelihood ratio. Crucially, this ratio is the exponential transform of the per-demonstration term in  $\Delta$  (Eq. 1). Therefore, sampling from the conditioned policy naturally favors high- $\Delta$  traces without any explicit calculation of  $\Delta$ . We formalize this relationship below.

**Theoretical Result.** Let  $R(q, r) \in \{0, 1\}$  be correctness reward. In-Context RLVR objective is

$$J_{\text{IC}}(\theta) = \mathbb{E}_q \mathbb{E}_{e \sim \mathcal{E}, r \sim \pi_{\theta}(\cdot|e, q)} [R(q, r)]. \quad (2)$$

We show (Appendix E) that this is equivalent to

$$J(\theta) = \mathbb{E}_q \mathbb{E}_{r \sim \pi_{\theta}(\cdot|q)} [R(q, r) \cdot w(q, r)], \quad (3)$$

where the weight factor satisfies  $\log w(q, r) \approx \Delta(q, r) + c$  for some model-specific constant  $c$ , which implies that  $w(q, r) \propto \exp(\Delta(q, r))$ .

The reweighted reward  $R \cdot w$  implies a two-stage selection mechanism. First, the binary reward  $R$  filters out traces with incorrect answers, assigning them zero rewards. Second, among traces getting correct answers, the weight  $w \propto \exp(\Delta)$  differentiates reasoning quality based on  $\Delta$ , assigning higher weights to high-quality traces and lower weights to low-quality ones. Through this equivalence, while the training is explicitly based on  $\pi_{\theta}(r|e, q)$ , In-Context RLVR implicitly optimizes the base policy  $\pi_{\theta}(\cdot|q)$  with rewards reweighted by  $\Delta$ .

Eq.3 shows that our method employs the model’s own ICL ability to guide optimization, with the policy serving as both the learner and the implicit quality evaluator. A natural concern is whether Evidence Gain remains a valid quality signal as the policy evolves, since our validation in Section 2 uses a

<sup>1</sup>Computing  $\Delta$  for  $\sim 12\text{K}$  samples over 100 demonstrations requires approximately 80 hours on H800.

fixed model. We address this in Section 4, showing that the correlation between Evidence Gain and reasoning quality remains stable throughout training.

Notably, while  $J_{\text{IC}}$  and  $J$  are equivalent in the sense of *expectation*, they differ in the sense of *variance*. The reweighting term  $w \propto \exp(\Delta)$  in  $J$  would introduce prohibitive reward variance during rollouts.  $J_{\text{IC}}$  avoids this potential training instability by shifting the sampling policy directly.

## 4 Experiments

To validate our framework, we combine In-Context RLVR with DAPO (Yu et al., 2025), yielding **IC-DAPO**. We choose DAPO as the backbone because it is a widely adopted RLVR method whose key techniques (e.g., clip-higher) have been adopted by many subsequent methods (Yue et al., 2025; Su et al., 2025b), making it a representative baseline for evaluating input-side modifications. All details of this section are provided in Appendix B.

### 4.1 Setup

**Dataset.** Our training data is derived from KlearReasoner-MathSub-30K (Su et al., 2025a), which contains 30K mathematical reasoning problems. We partition training data into three disjoint subsets: (1) a **training set** for policy optimization, (2) a **demonstration set**  $\mathcal{E}$  containing 1,082 question-reasoning pairs used for demonstration during IC-DAPO training, and (3) a **held-out set**  $\mathcal{E}_0$  of 100 additional examples reserved for the correlation analysis in §4.3. Both  $\mathcal{E}$  and  $\mathcal{E}_0$  are constructed following the procedure described in §2.

**Baselines.** We compare against several popular RLVR methods, including both standard outcome-based algorithms and more advanced objective-modifying variants. This selection directly tests whether our input-side modification can match algorithmic innovations in policy optimization. We exclude PRM-based methods as they require costly overhead that our method aims to avoid.

**Training and Evaluation.** We train DeepSeek-R1-Distill-Qwen at 1.5B and 7B scales. We conduct evaluations across various authoritative mathematical reasoning benchmarks, including AIME24, AIME25, HMMT25, MATH500 (Lightman et al., 2024), AMC23 and OlympiadBench (He et al., 2024). Following Su et al. (2025b), we report avg@4 scores on MATH500 and OlympiadBench, and avg@32 scores on all other benchmarks. Cru-

Method	AIME24	AIME25	HMMT25	MATH500	AMC23	Olympiad	Average	Time/Step (s)
<b>DS-R1-Distill-Qwen-1.5B</b>	29.2	24.1	13.1	86.0	73.7	51.8	46.3	–
+ GRPO (Shao et al., 2024)	33.4	28.1	16.6	88.3	79.3	56.2	50.3	457.4
+ DAPO (Yu et al., 2025)	40.0	28.4	19.2	90.0	84.4	61.6	53.9	459.6
+ CISPO (MiniMax et al., 2025)	32.9	25.1	13.2	85.8	80.9	54.9	48.8	466.3
+ GSPO (Zheng et al., 2025)	42.5	<u>33.6</u>	19.0	90.3	<u>85.9</u>	<b>62.6</b>	<u>55.7</u>	437.3
+ CE-GPPO (Su et al., 2025b)	<u>42.8</u>	32.5	<b>20.5</b>	<b>91.0</b>	85.8	61.8	<u>55.7</u>	464.0
+ IC-DAPO (Ours)	<b>45.6</b> ↑	<b>34.2</b> ↑	<u>19.7</u> ↑	<u>90.6</u> ↑	<b>86.2</b> ↑	<u>62.1</u> ↑	<b>56.4</b> ↑	477.2
<b>DS-R1-Distill-Qwen-7B</b>	54.5	39.1	26.2	93.6	90.6	67.0	61.8	–
+ GRPO (Shao et al., 2024)	55.3	40.3	24.5	93.7	88.8	65.6	61.4	305.6
+ DAPO (Yu et al., 2025)	62.0	45.9	27.4	94.1	92.3	69.9	65.3	303.1
+ CE-GPPO (Su et al., 2025b)	<u>64.2</u>	<b>50.3</b>	28.9	<u>95.3</u>	<u>93.3</u>	<u>71.6</u>	<u>67.3</u>	292.5
+ IC-DAPO (Ours)	<b>66.5</b> ↑	<u>49.8</u> ↑	<u>29.4</u> ↑	<b>95.6</b> ↑	<b>93.7</b> ↑	<b>71.7</b> ↑	<b>67.8</b> ↑	315.6

Table 1: Performance comparison across mathematical reasoning benchmarks. **Bold** and underline indicate the best and second-best results respectively. ↑ denotes improvement over the DAPO baseline. Notably, training times are incomparable across scales due to different GPU configurations (32 vs. 128 GPUs). Given this high GPU requirements at 7B scales, GSPO and CISPO are evaluated only at 1.5B.

cially, *all evaluation is conducted in zero-shot mode*, ensuring fair comparison with baselines.

## 4.2 Main Results

Table 1 presents benchmark performance. IC-DAPO outperforms DAPO by +2.5 average points at both scales, with gains particularly pronounced on competition benchmarks: +5.6 on AIME24 and +5.8 on AIME25 for the 1.5B model. This supports our hypothesis that implicit quality-based reweighting helps more on challenging problems where correct-but-low-quality traces are most harmful.

Beyond improvements over DAPO, IC-DAPO also matches or exceeds methods that modify the RL objective (e.g., GSPO, CISPO), achieving the highest average score at both scales while only altering the *input distribution*. This suggests that input-side modification constitutes an improvement axis orthogonal to policy optimization algorithms. We further compare wall-clock training time per training step and find that IC-DAPO incurs slight overhead (<5%), confirming its practicality.

## 4.3 Analysis

Our theory (§3) predicts that In-Context RLVR implicitly upweights high- $\Delta$  traces. To verify this, we track training dynamics by computing  $\Delta$  on the held-out set  $\mathcal{E}_0$  and assessing reasoning quality via Deepseek-V3.2, following procedures in §2.

**Q1: Does implicit reweighting occur?** Figure 2 (left) shows that mean Evidence Gain increases steadily under IC-DAPO throughout training, while DAPO exhibits smaller and slower growth. This confirms that the conditioned objective steers the policy toward traces with higher demonstration utility, exactly as predicted.

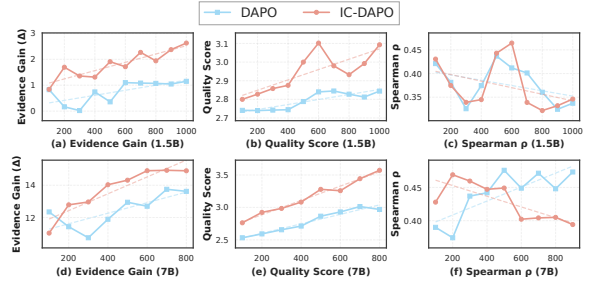


Figure 2: Training dynamics of Evidence Gain, quality score, and their correlation ( $\rho$ ) across training steps.

**Q2: Does this improve reasoning quality?** Figure 2 (middle) shows higher  $\Delta$  corresponds to improved quality scores. Note that  $\mathcal{E}_0$  used for evaluation is disjoint from  $\mathcal{E}$  used for training, ensuring unbiased evaluation. This shows that by upweighting traces with high teaching utility, we encourage better reasoning rather than merely correct answers.

**Q3: Is Evidence Gain valid throughout training?** Finally, we address the concern from §3. Figure 2 (right) shows that the Spearman correlation between  $\Delta$  and quality remains stable (around  $\rho \approx 0.4$ ) across training steps, confirming that the policy model’s intrinsic ICL signal remains a robust quality indicator as training progresses.

## 5 Conclusions

We introduce Evidence Gain, a quality signal that measures reasoning quality based on the policy model’s intrinsic ICL ability. To leverage this signal, we propose In-Context RLVR, which prepends demonstrations during training and implicitly reweights rewards by Evidence Gain to encourage high-quality traces. Experiments confirm improvements in both accuracy and reasoning quality over standard RLVR, providing a practical approach to improve reasoning quality in RLVR.

## 325 Limitations

326 This work has two main limitations. First, although  
327 In-Context RLVR demonstrates consistent improve-  
328 ments across mathematical reasoning benchmarks,  
329 its generalization to other reasoning-intensive do-  
330 mains such as STEM problem-solving remains an  
331 open question due to computational constraints.  
332 We focus on mathematical reasoning as it repre-  
333 sents a major research direction in LLM reinforce-  
334 ment learning, and we plan to pursue broader do-  
335 main evaluation in future work. Second, construct-  
336 ing the demonstration set requires access to a strong  
337 model (e.g., DeepSeek-R1) for generating high-  
338 quality reference traces. Alternative construction  
339 strategies that reduce this dependency should be  
340 further developed.

## 341 References

342 Yang Chen, Zhuolin Yang, Zihan Liu, Chankyu Lee,  
343 Peng Xu, Mohammad Shoeybi, Bryan Catanzaro,  
344 and Wei Ping. 2025. [Acereason-nemotron: Advancing math and code reasoning through reinforcement learning](#). *Preprint*, arXiv:2505.16400.

347 Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan,  
348 Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen Fan,  
349 Huayu Chen, Weize Chen, Zhiyuan Liu, Hao Peng,  
350 Lei Bai, Wanli Ouyang, Yu Cheng, Bowen Zhou, and  
351 Ning Ding. 2025. [The entropy mechanism of reinforcement learning for reasoning language models](#). *Preprint*, arXiv:2505.22617.

354 DeepSeek-AI. 2024. [Deepseek-v3 technical report](#).  
355 *Preprint*, arXiv:2412.19437.

356 DeepSeek-AI, Aixin Liu, Aoxue Mei, Bangcai Lin,  
357 Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao  
358 Wu, Bowei Zhang, Chaofan Lin, Chen Dong,  
359 Chengda Lu, Chenggang Zhao, Chengqi Deng, Chen-  
360 hao Xu, Chong Ruan, Damai Dai, Daya Guo, Dejian  
361 Yang, and 245 others. 2025. [Deepseek-v3.2: Pushing the frontier of open large language models](#). *Preprint*, arXiv:2512.02556.

364 Heejin Do, Jaehui Hwang, Dongyoon Han, Seong Joon  
365 Oh, and Sangdoon Yun. 2025. [What defines good reasoning in llms? dissecting reasoning steps with multi-aspect evaluation](#). *Preprint*, arXiv:2510.20603.

368 Olga Golovneva, Moya Chen, Spencer Poff, Martin  
369 Corredor, Luke Zettlemoyer, Maryam Fazel-Zarandi,  
370 and Asli Celikyilmaz. 2023. [ROSCOE: A suite of metrics for scoring step-by-step reasoning](#). In *ICLR 2023*. OpenReview.net.

373 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,  
374 Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang,  
375 Shirong Ma, Xiao Bi, Xiaokang Zhang, Xingkai Yu,

376 Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhu-  
377 oshu Li, Ziyi Gao, Aixin Liu, and 175 others. 2025a.  
378 Deepseek-r1 incentivizes reasoning in llms through  
379 reinforcement learning. *Nat.*, 645(8081):633–638.

380 Jiaxing Guo, Wenjie Yang, Shengzhong Zhang, Tong-  
381 shan Xu, Lun Du, Da Zheng, and Zengfeng Huang.  
382 2025b. [Right is not enough: The pitfalls of outcome supervision in training llms for math reasoning](#). *Preprint*, arXiv:2506.06877.

385 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu,  
386 Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie  
387 Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan  
388 Liu, and Maosong Sun. 2024. [OlympiadBench: A challenging benchmark for promoting AGI with olympiad-level bilingual multimodal scientific problems](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3828–3850.

394 Jujie He, Jiakai Liu, Chris Yuhao Liu, Rui Yan, Chaojie  
395 Wang, Peng Cheng, Xiaoyu Zhang, Fuxiang Zhang,  
396 Jiacheng Xu, Wei Shen, Siyuan Li, Liang Zeng, Tian-  
397 wen Wei, Cheng Cheng, Bo An, Yang Liu, and Yahui  
398 Zhou. 2025. [Skywork open reasoner 1 technical report](#). *Preprint*, arXiv:2505.22312.

400 Nathan Lambert, Jacob Morrison, Valentina Pyatkin,  
401 Shengyi Huang, Hamish Ivison, Faeze Brahman,  
402 Lester James V. Miranda, Alisa Liu, Nouha Dziri,  
403 Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf,  
404 Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras,  
405 and Oyvind Tafjord. 2025. [Tulu 3: Pushing frontiers in open language model post-training](#). *Preprint*, arXiv:2411.15124.

408 Dacheng Li, Shiyi Cao, Tyler Griggs, Shu Liu, Xi-  
409 angxi Mo, Eric Tang, Sumanth Hegde, Kourosh  
410 Hakhamaneshi, Shishir G. Patil, Matei Zaharia,  
411 Joseph E. Gonzalez, and Ion Stoica. 2025. [Llms can easily learn to reason from demonstrations structure, not content, is what matters!](#) *Preprint*, arXiv:2502.07374.

415 Jia LI, Edward Beeching, Lewis Tunstall, Ben Lipkin,  
416 Roman Soletskyi, Shengyi Costa Huang, Kashif Ras-  
417 sul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin,  
418 Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lam-  
419 ple, and Stanislas Polu. 2024. [Numinamath](#). <https://huggingface.co/AI-MO/NuminaMath-CoT>.  
420 Technical report available at [https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina\\_dataset.pdf](https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf).

424 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harri-  
425 son Edwards, Bowen Baker, Teddy Lee, Jan Leike,  
426 John Schulman, Ilya Sutskever, and Karl Cobbe.  
427 2024. [Let’s verify step by step](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.

430 Man Luo, Xin Xu, Yue Liu, Panupong Pasupat, and  
431 Mehran Kazemi. 2024. [In-context learning with retrieved demonstrations for language models: A survey](#). *Preprint*, arXiv:2401.11624.

434	Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi,	Shijie Xia, Xuefeng Li, Yixin Liu, Tongshuang Wu, and	490
435	William Y. Tang, Manan Roongta, Colin Cai, Jeffrey	Pengfei Liu. 2025. Evaluating mathematical reason-	491
436	Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica.	ing beyond accuracy. In <i>AAAI-25, Sponsored by the</i>	492
437	2025. Deepscaler: Surpassing o1-preview with a	<i>Association for the Advancement of Artificial Intelli-</i>	493
438	1.5b model by scaling rl. Notion Blog.	<i>gence, February 25 - March 4, 2025, Philadelphia,</i>	494
		<i>PA, USA, pages 27723–27730. AAAI Press.</i>	495
439	Monte MacDiarmid, Benjamin Wright, Jonathan Ue-	Wei Xiong, Hanning Zhang, Chenlu Ye, Lichang Chen,	496
440	sato, Joe Benton, Jon Kutasov, Sara Price, Naia Bous-	Nan Jiang, and Tong Zhang. 2025. Self-rewarding	497
441	cal, Sam Bowman, Trenton Bricken, Alex Cloud, and	correction for mathematical reasoning. <i>Preprint,</i>	498
442	Carson Denison. 2025. Natural emergent misalign-	arXiv:2502.19613.	499
443	ment from reward hacking in production rl. <i>Preprint,</i>		
444	arXiv:2511.18397.		
445	Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe,	An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao,	500
446	Mike Lewis, Hannaneh Hajishirzi, and Luke Zettle-	Bowen Yu, Chengpeng Li, Dayiheng Liu, Jian-	501
447	moyer. 2022. Rethinking the role of demonstrations:	hong Tu, Jingren Zhou, Junyang Lin, Keming Lu,	502
448	What makes in-context learning work? In <i>Proceed-</i>	Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang	503
449	<i>ings of the 2022 Conference on Empirical Methods</i>	Ren, and Zhenru Zhang. 2024. Qwen2.5-math tech-	504
450	<i>in Natural Language Processing.</i>	nical report: Toward mathematical expert model via	505
		self-improvement. <i>Preprint, arXiv:2409.12122.</i>	506
451	MiniMax, :, Aili Chen, Aonian Li, Bangwei Gong,	Chenlu Ye, Zhou Yu, Ziji Zhang, Hao Chen, Narayanan	507
452	Binyang Jiang, Bo Fei, Bo Yang, Boji Shan,	Sadagopan, Jing Huang, Tong Zhang, and Anurag	508
453	Changqing Yu, Chao Wang, Cheng Zhu, Chengjun	Beniwal. 2025. Beyond correctness: Harmonizing	509
454	Xiao, Chengyu Du, Chi Zhang, Chu Qiao, Chun-	process and outcome rewards through rl training.	510
455	hao Zhang, Chunhui Du, Congchao Guo, and 109	<i>Preprint, arXiv:2509.03403.</i>	511
456	others. 2025. Minimax-m1: Scaling test-time com-		
457	pute efficiently with lightning attention. <i>Preprint,</i>	Seonghyeon Ye, Doyoung Kim, Sungdong Kim, Hyeon-	512
458	arXiv:2506.13585.	bin Hwang, Seungone Kim, Yongrae Jo, James	513
		Thorne, Juho Kim, and Minjoon Seo. 2024. Flask:	514
459	Youssef Mroueh. 2025. Reinforcement learning	Fine-grained language model evaluation based on	515
460	with verifiable rewards: Grpo’s effective loss,	alignment skill sets. <i>Preprint, arXiv:2307.10928.</i>	516
461	dynamics, and success amplification. <i>Preprint,</i>		
462	arXiv:2503.06639.		
463	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan,	517
464	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	Xiaochen Zuo, Yu Yue, Weinan Dai, and Tiantian	518
465	Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024.	Fan. 2025. Dapo: An open-source llm rein-	519
466	Deepseekmath: Pushing the limits of mathemat-	forcement learning system at scale. <i>Preprint,</i>	520
467	ical reasoning in open language models. <i>Preprint,</i>	arXiv:2503.14476.	521
468	arXiv:2402.03300.		
469	Zhenpeng Su, Leiyu Pan, Xue Bai, Dening Liu,	Yu Yue, Yufeng Yuan, Qiyang Yu, Xiaochen Zuo, Ruofei	522
470	Guanting Dong, Jiaming Huang, Wenping Hu,	Zhu, Wenyuan Xu, Jiase Chen, Chengyi Wang,	523
471	Fuzheng Zhang, Kun Gai, and Guorui Zhou. 2025a.	TianTian Fan, Zhengyin Du, Xiangpeng Wei, Xi-	524
472	Klear-reasoner: Advancing reasoning capability	angyu Yu, Gaohong Liu, Juncai Liu, Lingjun Liu,	525
473	via gradient-preserving clipping policy optimization.	Haibin Lin, Zhiqi Lin, Bole Ma, Chi Zhang, and 8	526
474	<i>Preprint, arXiv:2508.07629.</i>	others. 2025. Vapo: Efficient and reliable reinforce-	527
		ment learning for advanced reasoning tasks. <i>Preprint,</i>	528
475	Zhenpeng Su, Leiyu Pan, Minxuan Lv, Yuntao Li, Wen-	arXiv:2504.05118.	529
476	ping Hu, Fuzheng Zhang, Kun Gai, and Guorui Zhou.		
477	2025b. Ce-gppo: Coordinating entropy via gradient-	Kongcheng Zhang, Qi Yao, Shunyu Liu, Yingjie Wang,	530
478	preserving clipping policy optimization in reinforce-	Baisheng Lai, Jieping Ye, Mingli Song, and Dacheng	531
479	ment learning. <i>Preprint, arXiv:2509.20712.</i>	Tao. 2025a. Consistent paths lead to truth: Self-	532
		rewarding reinforcement learning for llm reasoning.	533
480	Zhenpeng Su, Leiyu Pan, Minxuan Lv, Tiehua Mei, Zi-	In <i>Advances in Neural Information Processing Sys-</i>	534
481	jia Lin, Yuntao Li, Wenping Hu, Ruiming Tang, Kun	<i>tems.</i>	535
482	Gai, and Guorui Zhou. 2025c. Entropy ratio clipping		
483	as a soft global constraint for stable reinforcement	Zhenru Zhang, Chujie Zheng, Yangzhen Wu, Beichen	536
484	learning. <i>Preprint, arXiv:2512.05591.</i>	Zhang, Runji Lin, Bowen Yu, Dayiheng Liu, Jin-	537
		gren Zhou, and Junyang Lin. 2025b. The lessons of	538
485	Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai	developing process reward models in mathematical	539
486	Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang Sui.	reasoning. <i>Preprint, arXiv:2501.07301.</i>	540
487	2024. Math-shepherd: Verify and reinforce llms		
488	step-by-step without human annotations. <i>Preprint,</i>	Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui	541
489	arXiv:2312.08935.	Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong	542
		Liu, Rui Men, An Yang, Jingren Zhou, and Jun-	543
		yang Lin. 2025. Group sequence policy optimization.	544
		<i>Preprint, arXiv:2507.18071.</i>	545

## A Related Work

**Reinforcement Learning with Verifiable Rewards.** RLVR has become a dominant paradigm for improving LLM reasoning (Lambert et al., 2025; Guo et al., 2025a). By using rule-based correctness signals, RLVR avoids reward hacking and scales well. GRPO (Shao et al., 2024) optimizes policies via group-based reward normalization without requiring a separate critic model. Subsequent work addresses entropy collapse (Yu et al., 2025), sequence-level optimization (Zheng et al., 2025), and training stability (MiniMax et al., 2025). Despite the advances, a fundamental limitation persists. Binary correctness rewards assign equal reward to correct solutions regardless of reasoning quality, potentially reinforcing spurious traces that get correct answers through flawed logic (Zhang et al., 2025b; MacDiarmid et al., 2025).

**Process Reward Models.** PRMs address this limitation by providing step-level feedback. Lightman et al. (2024) demonstrate substantial gains from process supervision over outcome supervision, though their approach required approximately 800,000 human-annotated step labels. Automated alternatives such as Math-Shepherd (Wang et al., 2024) construct step-level labels through Monte Carlo estimation, trading annotation cost for computational overhead from repeated rollouts. These methods share a common constraint. Obtaining reliable process signals demands either substantial human effort or significant compute resources.

**Self-Derived Reasoning Quality Signals.** An emerging direction leverages signals derived from the model itself rather than external supervision. Zhang et al. (2025a) construct intrinsic rewards from trajectory consistency and volatility, which requires computing distances between each intermediate state and all distinct final answers across sampled rollouts, incurring  $O(NTK)$  additional forward passes per prompt beyond standard algorithm. Xiong et al. (2025) train models to perform iterative self-correction loops, where the model sequentially detects errors, revises outputs, and decides when to terminate. In contrast, our method measures reasoning quality through demonstration utility, motivated by findings that high-quality reasoning traces can serve as effective in-context demonstrations (Min et al., 2022; Li et al., 2025). Crucially, this signal can be integrated implicitly into the training objective via In-Context RLVR,

adding less than 5% overhead without any explicit quality computation. Importantly, our approach differs fundamentally from demonstration selection methods (Luo et al., 2024), which develop retrieval or optimization strategies to identify the best demonstrations for each input query. In contrast, we leverage demonstrations to assess the quality of candidate reasoning traces; since Evidence Gain (Eq.1) is defined as an average over the validation set, all demonstrations are equally important for every query in our method.

## B Experimental Details

### B.1 Datasets and Preprocessing

Our training data is derived from KlearReasoner-MathSub-30K (Su et al., 2025a), which contains approximately 30K high-quality mathematical reasoning problems collected from several curated sources, including Skywork-OR1 (He et al., 2025), Acereason (Chen et al., 2025), NuminaMath (Li et al., 2024), and DeepScaleR (Luo et al., 2025). To mitigate potential data contamination, the dataset has been processed with 9-gram deduplication against the evaluation benchmarks.

We first randomly partition 28k samples from the full dataset for policy optimization. For the remaining 2k samples, we generate reasoning traces using DeepSeek-R1-0528 (Guo et al., 2025a) and filter the outputs using rule-based validators to retain only those with correct final answers, yielding approximately 1,200 valid examples. From this filtered set, we randomly select 100 examples to form the held-out set  $\mathcal{E}_0$  for the correlation analysis in Section 4.3, with the remaining 1,082 examples forming the demonstration set  $\mathcal{E}$  used for prepending demonstrations during training. For each example in  $\mathcal{E}$  and  $\mathcal{E}_0$ , we extract the content following `</think>` as the reference reasoning trace following the procedure described in Section 2. We manually verify the quality of these solutions.

### B.2 Details for Main Experiments

We train DeepSeek-R1-Distill-Qwen-1.5B<sup>2</sup> and DeepSeek-R1-Distill-Qwen-7B<sup>3</sup>. We evaluate on six mathematical reasoning benchmarks: AIME24, AIME25, HMMT25, MATH500, AMC23, and OlympiadBench. For evaluation metrics, we report avg@4 scores on MATH500 (Lightman et al.,

<sup>2</sup><https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-1.5B>

<sup>3</sup><https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-7B>

2024) and OlympiadBench (He et al., 2024), and avg@32 scores on all other benchmarks, following prior work (Su et al., 2025b). At inference, we set the maximum generation length to 32k tokens for AIME24 and AIME25, and 16k tokens for the other datasets. For answer extraction, we follow the standard practice adopted in Yang et al. (2024): parsing the contents enclosed within the `\boxed{\}` structure in model outputs to identify the final answer. Answer correctness is judged by math-verify<sup>4</sup>, which performs symbolic comparison to handle equivalent mathematical expressions.

For checkpoint selection, we save checkpoints every 10 training steps and select the checkpoint with the best performance on AIME25 to report results on all other benchmarks. All evaluation is conducted in zero-shot mode without any demonstrations, ensuring fair comparison with baseline methods and validating that our approach requires no modification to deployment. This zero-shot evaluation also empirically confirms our theoretical claim in Section 3: while demonstrations are used during training to enable implicit quality reweighting, In-Context RLVR implicitly optimizes the base policy which can operate without demonstrations.

### B.3 Details for Training Dynamics Analysis

Theoretical analysis in Section 3 shows that In-Context RLVR implicitly reweights rewards by Evidence Gain, with the reweighted reward  $R \cdot w$  implying a two-stage selection mechanism: first, the binary reward  $R$  filters out traces with incorrect answers ( $R = 0$ ); second, among correct traces ( $R = 1$ ), the weight  $w \propto \exp(\Delta)$  differentiates reasoning quality. Therefore, Evidence Gain is designed to distinguish reasoning quality among correct solutions. To validate this implicit reweighting mechanism empirically (Section 4.3), we track both Evidence Gain and reasoning quality scores exclusively on traces with correct final answers.

Specifically, we use checkpoints from both DAPO and IC-DAPO at 1.5B and 7B scales. Every 100 training steps, we randomly sample 100 queries from the training set. Importantly, at each chosen step, DAPO and IC-DAPO share the same set of sampled queries to ensure fair comparison. For each query, we generate 8 rollouts and retain only those with correct final answers. We then compute Evidence Gain on the held-out set  $\mathcal{E}_0$ , and assess reasoning quality using DeepSeek-V3.2 as

described in Appendix B.4. Finally, we compute the Spearman correlation  $\rho$  between Evidence Gain and quality scores to verify that Evidence Gain remains a valid quality signal throughout training.

### B.4 LLM-based Quality Evaluation

To automatically assess reasoning quality, we employ DeepSeek-V3.2 (DeepSeek-AI et al., 2025) as an LLM-based evaluator. Our evaluation rubric is informed by prior work on solution quality assessment (Ye et al., 2024; Xia et al., 2025; Golovneva et al., 2023). To ensure comprehensive coverage, we define eight complementary dimensions as shown in Table 2. For each reasoning trace, the evaluator assigns a score from 1 to 5 on each dimension along with explicit textual explanations for justification, and finally provides an overall quality score from 1 to 5. The complete prompt template is provided in Appendix D.

Dimension	Definition
Repetition	Same steps or ideas repeated
Redundancy	Unnecessary or verbose content
Logical Consistency	Contradictions
Relevance	Off-topic exploration
CoT-Ans Alignment	Answer derived from reasoning
Reasoning Rigor	Claims justified without leaps
Clarity	Easy to follow, well-structured
Completeness	All necessary steps present

Table 2: Quality evaluation dimensions.

### B.5 Baseline Implementation

**GRPO** optimizes policies via group-based reward normalization without requiring a separate critic model. Following Shao et al. (2024), we adopt symmetric clipping bounds with  $\epsilon = 0.2$ .

**DAPO** extends GRPO by introducing asymmetric clipping bounds and dynamic sample filtering to mitigate entropy collapse. Following Yu et al. (2025), we set the lower and upper clipping thresholds to  $\epsilon_{\text{low}} = 0.2$  and  $\epsilon_{\text{high}} = 0.28$ , respectively.

**CISPO** applies clipping directly to the importance sampling weights rather than to the final policy update. Following Cui et al. (2025), we set symmetric clipping bounds with  $\epsilon = 0.2$ .

**GSPO** employs a sequence-level importance ratio to enhance training stability and scalability. Following Zheng et al. (2025), we set the lower and upper clipping thresholds to  $\epsilon_{\text{low}} = 0.0003$  and  $\epsilon_{\text{high}} = 0.0004$ , respectively.

<sup>4</sup><https://github.com/huggingface/Math-Verify>

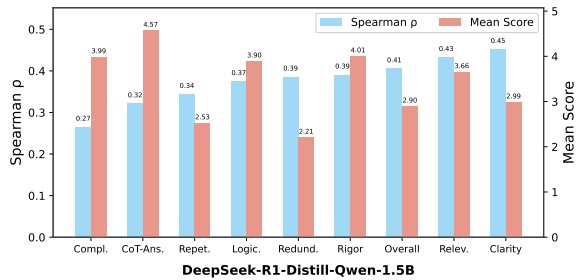


Figure 3: Spearman correlation between Evidence Gain and individual quality dimensions, alongside mean scores for each dimension. Evaluated on DeepSeek-R1-Distill-Qwen-1.5B with the same setup as Section 2.

**CE-GPPO** reintroduces gradient signals from tokens outside the clipping interval in a bounded manner through a stop-gradient operation, enabling fine-grained control over policy entropy dynamics. We directly report results from the original paper (Su et al., 2025b). As this work provides two sets of evaluation results under different configurations ( $\beta_1 = 0.5, \beta_2 = 1$  and  $\beta_1 = 0.75, \beta_2 = 1$ ), we report their average scores in Table 1.

## C Additional Experiments

### C.1 What Does Evidence Gain Capture?

To investigate which aspects of reasoning quality Evidence Gain most effectively captures, we conduct a fine-grained correlation analysis across the eight dimensions defined in our evaluation rubric (Appendix B.4). Figure 3 presents results on DeepSeek-R1-Distill-Qwen-1.5B. We observe that all correlations are positive, ranging from  $\rho = 0.27$  to  $\rho = 0.45$ , suggesting that Evidence Gain reflects multiple aspects of reasoning quality rather than a single dimension.

Notably, Clarity ( $\rho = 0.45$ ) and Relevance ( $\rho = 0.43$ ) show stronger correlations than Completeness ( $\rho = 0.27$ ) and CoT-Answer Alignment ( $\rho = 0.32$ ). This pattern can be explained by our setup. Since we only evaluate traces that arrive at correct answers (Appendix B), these traces are already filtered for answer validity. As shown by the mean scores in Figure 3, Completeness (mean 3.99) and CoT-Answer Alignment (mean 4.57) exhibit high scores with limited variance among correct traces, because reaching the right answer typically requires including necessary steps and properly deriving the conclusion. In contrast, Clarity (mean 2.99) and Relevance (mean 3.66) show lower mean scores even among correct solutions, since a trace

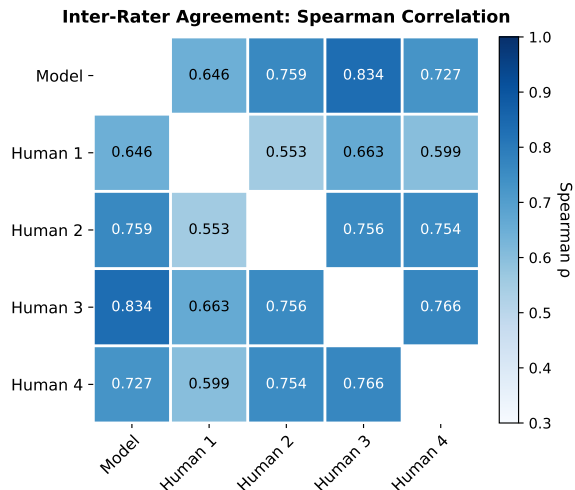


Figure 4: Spearman correlation matrix between DeepSeek-V3.2 quality scores and four human expert ratings on 100 sampled reasoning traces, together with inter rater correlations among experts.

can reach the right answer while still being poorly organized. This difference in variance explains why Evidence Gain shows stronger correlations with dimensions that have greater room to discriminate among correct traces.

### C.2 Human Evaluation

We further validate whether the automatic quality scores produced by DeepSeek-V3.2 are consistent with human judgement. We first randomly sample 100 question-reasoning pairs generated by Deepseek-R1-Distill-Qwen-1.5B whose final answers are verified correct by rule. We then ask **four** human experts to separately assign an overall quality score between 1 and 5 for each reasoning trace. Annotators follow the same rubric for Deepseek-V3.2 (Appendix B.4). Finally, we compute Spearman correlation coefficients between the DeepSeek-V3.2 scores and each human expert, as well as the correlations among the human experts.

Figure 4 shows strong agreement between DeepSeek-V3.2 and human experts. The correlations between DeepSeek-V3.2 and individual experts fall between 0.65 and 0.83, with an average of 0.74, which is comparable to human expert agreement. Overall, DeepSeek-V3.2 performs within the variability of human judgement on this task, supporting its reliability as an automatic evaluator in our experiments. We attribute this robustness to the comprehensive and fine-grained quality rubric, which decomposes reasoning quality into multiple distinguishable dimensions and thus enables more

Method	AIME24	AIME25	HMMT25	MATH500	AMC23	Olympiad	Average
<b>DS-R1-Distill-Qwen-1.5B</b>	29.2	24.1	13.1	86.0	73.7	51.8	46.3
+ DAPO (Yu et al., 2025)	40.0	28.4	19.2	90.0	84.4	61.6	53.9
+ IC-DAPO (V3.1)	<u>44.5</u> ↑	<u>32.3</u> ↑	<u>19.5</u> ↑	<u>90.3</u> ↑	<u>85.8</u> ↑	<u>61.7</u> ↑	<u>55.7</u> ↑
+ IC-DAPO (R1)	<b>45.6</b> ↑	<b>34.2</b> ↑	<b>19.7</b> ↑	<b>90.6</b> ↑	<b>86.2</b> ↑	<b>62.1</b> ↑	<b>56.4</b> ↑
<b>DS-R1-Distill-Qwen-7B</b>	54.5	39.1	26.2	93.6	90.6	67.0	61.8
+ DAPO (Yu et al., 2025)	62.0	45.9	27.4	94.1	92.3	69.9	65.3
+ IC-DAPO (V3.1)	<u>63.3</u> ↑	<u>47.5</u> ↑	<u>29.2</u> ↑	<u>95.5</u> ↑	<u>92.6</u> ↑	<u>70.8</u> ↑	<u>66.4</u> ↑
+ IC-DAPO (R1)	<b>66.5</b> ↑	<b>49.8</b> ↑	<b>29.4</b> ↑	<b>95.6</b> ↑	<b>93.7</b> ↑	<b>71.7</b> ↑	<b>67.8</b> ↑

Table 3: Ablation on demonstration quality. IC-DAPO (R1) uses refined reasoning traces from DeepSeek-R1, while IC-DAPO (V3.1) uses solutions from DeepSeek-V3.1, a non-reasoning model. Both variants outperform DAPO, with R1 demonstrations yielding stronger results. **Bold** and underline indicate best and second-best results. ↑ denotes improvement over DAPO.

consistent judgements across different raters.

### C.3 Ablation Study on Demonstration Quality

Our main experiments in Section 4 construct the demonstration set  $\mathcal{E}$  using the refined content following DeepSeek-R1’s `</think>` tag. To investigate whether demonstration quality affects training outcomes, we construct an alternative set using DeepSeek-V3.1 (DeepSeek-AI, 2024), a strong non-reasoning model, for comparison. Without the reasoning capabilities that produce R1’s refined traces, V3.1 is expected to generate solutions of lower quality, making it a suitable baseline for studying the effect of demonstration quality. To ensure fair comparison, we use the same 1,082 questions from the demonstration set  $\mathcal{E}$  described in Section 4 and generate solutions using DeepSeek-V3.1. We generate multiple responses per question to ensure each question obtains a correct solution.

Table 3 shows that both IC-DAPO variants outperform DAPO, confirming In-Context RLVR’s robustness over different demonstration sources. Crucially, IC-DAPO (R1) always surpasses IC-DAPO (V3.1), with gains of +1.1 on AIME24 and +1.9 on AIME25 at 1.5B scale. From a theoretical perspective, since the reference solutions in validation set are used to compute Evidence Gain (Eq. 1), the quality of these references may affect the accuracy of  $\Delta$  as a quality signal, which in turn affects the effectiveness of implicit reweighting (Eq. 3).

### C.4 Training Dynamics

We present extended training dynamics for both 1.5B and 7B models in Figure 5. Across both model scales, IC-DAPO variants (i.e., IC-DAPO (R1) and IC-DAPO (V3.1)) consistently achieve higher accuracy on AIME24 and AIME25 throughout training while maintaining comparable entropy trajectories to the DAPO baseline. Notably, IC-DAPO

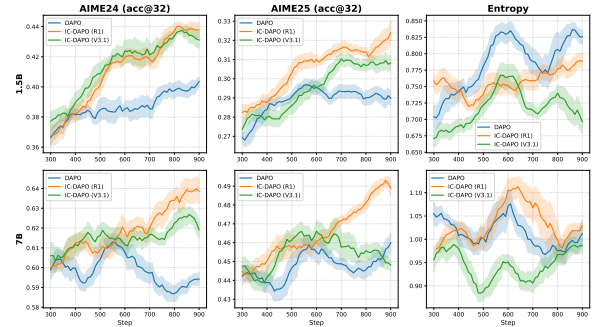


Figure 5: Training dynamics across 1.5B and 7B models. IC-DAPO variants consistently outperform DAPO on AIME24 and AIME25 while maintaining stable entropy throughout training.

(R1) demonstrates a clear advantage over standard DAPO, with the performance gap widening as training progresses. Comparing IC-DAPO (R1) and IC-DAPO (V3.1), we observe that R1-generated demonstrations yield better final performance, consistent with our findings in Section C.3 that higher-quality demonstrations lead to improved training outcomes. The entropy curves remain stable across all methods, indicating that In-Context RLVR does not compromise training stability while achieving superior accuracy, which aligns with the stability analysis in prior work (Yu et al., 2025).

## D Quality Evaluation Prompt

We present the prompt template used for reasoning quality evaluation with DeepSeek-V3.2 in Figure 6. The prompt assess reasoning traces across eight dimensions, including Repetition, Redundancy, Logical Consistency, Relevance, CoT-Answer Alignment, Reasoning Rigor, Clarity, and Completeness, with each scoring on a 1–5 scale. The template uses placeholders `{question}` and `{response}` which are populated with the specific math problem and corresponding reasoning trace during evaluation.

```

# Mathematical Reasoning Quality Evaluation

## Task
Evaluate the quality of mathematical reasoning (answer already verified correct).
Assess whether this reasoning would be valuable as a learning reference.

Key insight: Correct answer  $\neq$  good reasoning. Watch for warning signs:
- Frequent "wait", "hold on", "let me try again"
- Same calculation repeated multiple times
- Long explorations that don't contribute to the answer
- Answer appearing without clear derivation

Be strict. Reasoning that wanders, second-guesses itself, or reaches the right answer through
messy exploration is not high-quality, regardless of correctness.

## Dimensions (1-5 each)

| Dimension          | Assess                                                                 |
|-----|-----|
| Repetition         | Same steps or ideas repeated?                                       |
| Redundancy         | Unnecessary or verbose content?                                       |
| Logical Consistency | Contradictions or backtracks?                                       |
| Relevance          | Off-topic content or dead-end exploration?                           |
| CoT-Answer Alignment | Answer clearly derived from reasoning?                               |
| Reasoning Rigor    | All claims justified without leaps?                                    |
| Clarity            | Easy to follow and well-structured?                                   |
| Completeness       | All necessary steps present?                                         |

## Scoring

| Score | Meaning                                                                 |
|-----|-----|
| 5     | Excellent: Perfect textbook quality                                   |
| 4     | Good: Minor flaws, suitable as reference                             |
| 3     | Average: Clear flaws, but followable                                 |
| 2     | Weak: Major issues hurting pedagogical value                          |
| 1     | Poor: Guessed answer, chaotic flow                                   |

## Evaluate This

### Math Problem:
{question}

### Reasoning Process (Answer Verified Correct):
{response}

---

Output JSON only:
{
  "dimensions": {
    "repetition": {"score": <1-5>, "comment": "<evidence>"},
    "redundancy": {"score": <1-5>, "comment": "<evidence>"},
    "logical_consistency": {"score": <1-5>, "comment": "<evidence>"},
    "relevance": {"score": <1-5>, "comment": "<evidence>"},
    "cot_answer_alignment": {"score": <1-5>, "comment": "<evidence>"},
    "reasoning_rigor": {"score": <1-5>, "comment": "<evidence>"},
    "clarity": {"score": <1-5>, "comment": "<evidence>"},
    "completeness": {"score": <1-5>, "comment": "<evidence>"},
  },
  "overall_analysis": "<key findings, limiting factors>",
  "score": <1-5>
}

```

Figure 6: Quality Evaluation Prompt Template. Placeholders `{question}` and `{response}` are replaced with the actual math problem and reasoning trace during evaluation.

## E Proofs of Main Results

[This section of appendix uses single-column format for mathematical readability.]

This appendix provides complete derivations for the theoretical results in Section 3. We first establish notation, then derive the reweighting identity and characterize its relationship to Evidence Gain.

### E.1 Notations

We formalize the training setup. Let  $q \sim \mathcal{D}$  denote a question from the training distribution. Let  $\mathcal{E} = \{e^{(i)}\}_{i=1}^{|\mathcal{E}|}$  be the held-out validation set, where each demonstration  $e = (e_q, e_r)$  consists of a question  $e_q$  and a high-quality reference reasoning trace  $e_r$ . During In-Context RLVR training, a demonstration  $e$  is sampled from  $\mathcal{E}$  and prepended to  $q$ , after which the model generates a reasoning trace  $r \sim \pi_\theta(\cdot|e, q)$ .

### E.2 Bayesian Identity

We establish the key identity relating the conditioned policy to the base policy. The derivation relies on the following assumption, which reflects the independent sampling structure in our data construction.

**Assumption E.1.** *Providing only the demonstration question  $e_q$ , without its reasoning trace  $e_r$ , does not alter the distribution over reasoning traces for a training question  $q$ . Conversely, providing only the training question  $q$ , without any reasoning trace, does not alter the distribution over reasoning traces for the demonstration question  $e_q$ . Formally:*

$$\pi_\theta(r|e_q, q) = \pi_\theta(r|q), \quad (\text{A1})$$

$$\pi_\theta(e_r|q, e_q) = \pi_\theta(e_r|e_q). \quad (\text{A2})$$

**Remark.** This assumption is natural given the independent sampling of demonstrations and training questions. To see why (A1) holds, suppose  $e_q$  is “Solve  $x^2 - 5x + 6 = 0$ ” and  $q$  is “Compute  $\int \sin x dx$ .” The bare statement of  $e_q$  carries no information about integration techniques; it only indicates that the context involves math. Crucially, the model already knows this from observing  $q$  itself. Thus, conditioning on  $e_q$  alone provides no additional signal for solving  $q$ . A symmetric argument establishes (A2). While edge cases may exist where two questions happen to share methodological structure, so that  $e_q$  could bias the preferred style or method for solving  $q$ , such coincidences are rare under independent sampling and average out at scale. Thus (A1) and (A2) hold as statistical approximations.

**Empirical Validation.** We validate Assumption E.1 using DeepSeek-R1-Distill-Qwen-1.5B. We randomly select 100 question-reasoning pairs  $(q, r)$  from rollouts generated during training, and independently sample 100 additional questions  $\{q'\}$  from the dataset. We prepend each  $q'$  to each  $(q, r)$  resulting in 10,000 samples  $(S)$ . We measure the relative change in reasoning log-probability when prepending an randomly select question  $q'$ :

$$\delta = \frac{1}{|S|} \sum_{(q,r,q') \in S} \frac{|\log \pi_\theta(r|q', q) - \log \pi_\theta(r|q)|}{|\log \pi_\theta(r|q)|}. \quad (4)$$

Notably, we obtain  $\delta = 0.0384 < 5\%$ , confirming that independently sampled question statements have negligible influence on reasoning distributions on average.

Although Assumption E.1 suggests that question statements alone provide negligible cross example influence, a complete demonstration  $(e_q, e_r)$  does provide transferable information: the reasoning trace  $e_r$  may exhibit problem-solving patterns (e.g., algebraic manipulation, problem decomposition) that generalize across problems. This distinction is precisely what Evidence Gain captures, and it explains why the policy model’s ICL ability can serve as an effective quality signal.

**Lemma E.2 (Bayesian Identity).** *Under Assumption E.1, the conditioned policy admits the decomposition:*

$$\pi_\theta(r|e, q) = \pi_\theta(r|q) \cdot \frac{\pi_\theta(e_r|q, r, e_q)}{\pi_\theta(e_r|e_q)}. \quad (5)$$

*Proof.* We begin by expanding the conditioned policy  $\pi_\theta(r|e_q, e_r, q)$  using the definition of conditional probability. By Bayes' Rule, we have:

$$\pi_\theta(r|e_q, e_r, q) = \frac{\pi_\theta(r|e_q, q) \pi_\theta(e_r|q, r, e_q)}{\pi_\theta(e_r|q, e_q)}. \quad (6)$$

We now apply Assumption E.1 to simplify the right-hand side. By (A1), the numerator term  $\pi_\theta(r|e_q, q)$  reduces to  $\pi_\theta(r|q)$ , since observing  $e_q$  alone provides no additional information for generating  $r$ . By (A2), the denominator term  $\pi_\theta(e_r|q, e_q)$  reduces to  $\pi_\theta(e_r|e_q)$ , since observing  $q$  alone provides no additional information for generating  $e_r$ . Substituting these simplifications into Eq. (6), we obtain:

$$\pi_\theta(r|e_q, e_r, q) = \pi_\theta(r|q) \cdot \frac{\pi_\theta(e_r|q, r, e_q)}{\pi_\theta(e_r|e_q)}. \quad (7)$$

This completes the proof.  $\square$

### E.3 Reweighting Equivalence

We now establish that In-Context RLVR implicitly performs reward reweighting, showing how the policy model's ICL mechanism naturally upweights high-quality reasoning traces.

**Theorem E.3** (Implicit Reweighting). *The In-Context RLVR objective*

$$J_{\text{IC}}(\theta) = \mathbb{E}_q \mathbb{E}_{e \sim \mathcal{E}, r \sim \pi_\theta(\cdot|e, q)} [R(q, r)] \quad (8)$$

can be equivalently written as

$$J_{\text{IC}}(\theta) = \mathbb{E}_q \mathbb{E}_{r \sim \pi_\theta(\cdot|q)} [R(q, r) \cdot w(q, r)], \quad (9)$$

where the weight factor  $w(q, r) = \mathbb{E}_e[\exp(\Delta_e)]$  with  $\Delta_e = \log \pi_\theta(e_r|q, r, e_q) - \log \pi_\theta(e_r|e_q)$ .

Furthermore,  $w(q, r)$  and the Evidence Gain  $\Delta(q, r) = \mathbb{E}_e[\Delta_e]$  satisfy:

1. (**Lower bound**) By Jensen's inequality:  $w(q, r) \geq \exp(\Delta(q, r))$ .

2. (**Refined bound**)  $\log w(q, r) = \Delta(q, r) + \log(1 + \frac{1}{2} \text{Var}_e[\Delta_e]) + o(\text{Var}_e[\Delta_e])$ .

*Proof. Part I: Deriving the weight factor.* Invoking Lemma E.2 and assuming uniform sampling over  $\mathcal{E}$ , we expand:

$$J_{\text{IC}}(\theta) = \mathbb{E}_q \left[ \frac{1}{|\mathcal{E}|} \sum_e \sum_r \pi_\theta(r|e, q) \cdot R(q, r) \right] \quad (10)$$

$$= \mathbb{E}_q \left[ \frac{1}{|\mathcal{E}|} \sum_e \sum_r \pi_\theta(r|q) \cdot \frac{\pi_\theta(e_r|q, r, e_q)}{\pi_\theta(e_r|e_q)} \cdot R(q, r) \right] \quad (11)$$

$$= \mathbb{E}_q \left[ \sum_r \pi_\theta(r|q) \cdot R(q, r) \cdot \underbrace{\frac{1}{|\mathcal{E}|} \sum_e \frac{\pi_\theta(e_r|q, r, e_q)}{\pi_\theta(e_r|e_q)}}_{w(q, r)} \right]. \quad (12)$$

Writing  $\Delta_e = \log \pi_\theta(e_r|q, r, e_q) - \log \pi_\theta(e_r|e_q)$ , the weight becomes  $w(q, r) = \mathbb{E}_e[\exp(\Delta_e)]$ .

**Part II: Relating  $w(q, r)$  to  $\Delta(q, r)$ .** Since  $\exp(\cdot)$  is convex, Jensen's inequality gives  $\mathbb{E}_e[\exp(\Delta_e)] \geq \exp(\mathbb{E}_e[\Delta_e]) = \exp(\Delta(q, r))$ , establishing (i).

For (ii), we expand  $\exp(\Delta_e)$  around  $\Delta := \Delta(q, r)$  via Taylor series:

$$\exp(\Delta_e) = \exp(\Delta) \cdot \exp(\Delta_e - \Delta) \quad (13)$$

$$= \exp(\Delta) \cdot \left( 1 + (\Delta_e - \Delta) + \frac{1}{2}(\Delta_e - \Delta)^2 + o((\Delta_e - \Delta)^2) \right). \quad (14)$$

Taking expectations and using  $\mathbb{E}_e[\Delta_e - \Delta] = 0$ :

$$w(q, r) = \exp(\Delta) \cdot \left( 1 + \frac{1}{2} \text{Var}_e[\Delta_e] + o(\text{Var}_e[\Delta_e]) \right). \quad (15)$$

Taking logarithms on both sides and applying  $\log(1 + x) = x + o(x)$  yields (ii).  $\square$

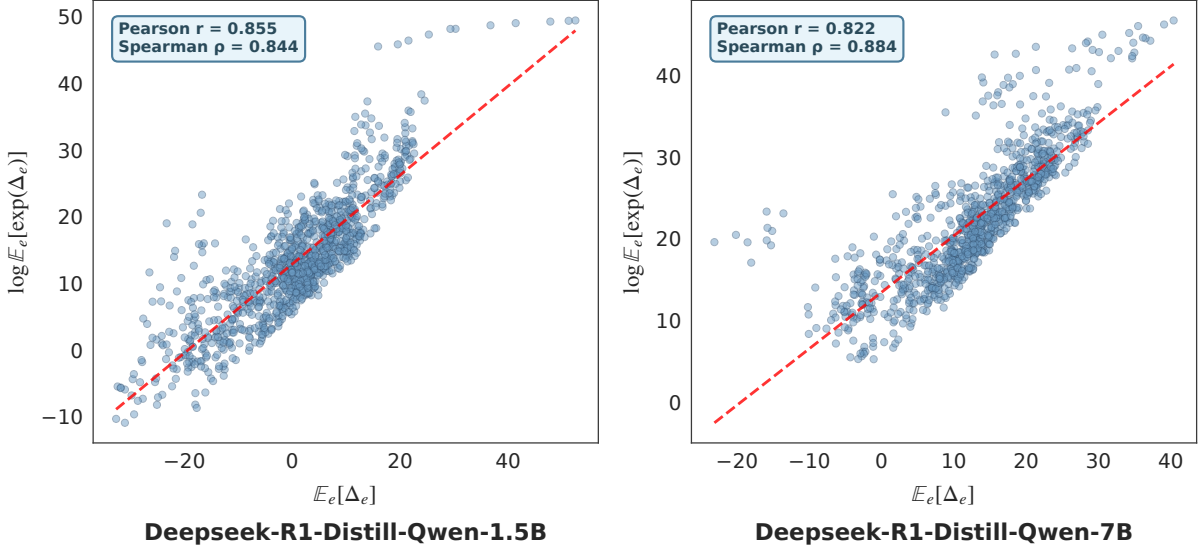


Figure 7: Empirical verification of the log-linear relationship between  $\log w(q, r) = \log \mathbb{E}_e[\exp(\Delta_e)]$  and  $\Delta(q, r) = \mathbb{E}_e[\Delta_e]$ . Strong Pearson correlations ( $r = 0.855$  for 1.5B,  $r = 0.822$  for 7B) confirm that Evidence Gain serves as a reliable proxy for the implicit weight.

**Interpretation.** Result (i) shows that  $\exp(\Delta(q, r))$  serves as a *lower bound* for the implicit weight  $w(q, r)$ , but does not quantify how tight this bound is. Result (ii) refines this by showing that the gap is controlled by  $\frac{1}{2}\text{Var}_e[\Delta_e]$ . Specifically,  $\log w(q, r)$  exceeds  $\Delta(q, r)$  by approximately  $\log(1 + \frac{1}{2}\text{Var}_e[\Delta_e]) > 0$ , a strictly positive correction.

**Empirical Observation.** Although Theorem E.3(ii) indicates that  $\log w(q, r)$  and  $\Delta(q, r)$  differ by a variance-dependent term, if  $\text{Var}_e[\Delta_e]$  remains relatively stable across different  $(q, r)$  pairs, the relationship simplifies to an approximate linear correspondence. To verify this, we conduct experiments using rollouts generated by DeepSeek-R1-Distill-Qwen at 1.5B and 7B (Guo et al., 2025a). Specifically, we randomly sample 1,100  $(q, r)$  pairs for the 1.5B model and 1,000 pairs for the 7B model. For each  $(q, r)$  pair, we compute  $\Delta(q, r) = \mathbb{E}_e[\Delta_e]$  and  $\log w(q, r) = \log \mathbb{E}_e[\exp(\Delta_e)]$ , then measure their correlation.

Figure 7 presents the results. We observe strong linear relationships. The 1.5B model yields Pearson  $r = 0.855$  (Spearman  $\rho = 0.844$ ), while the 7B model achieves  $r = 0.822$  (Spearman  $\rho = 0.884$ ). These high correlations confirm that the variance term contributes a near-constant offset across  $(q, r)$  pairs, validating the following practical characterization.

**Theorem E.4** (Reweighting Equivalence). *Under the conditions of Theorem E.3, the In-Context RLVR objective satisfies*

$$J_{\text{IC}}(\theta) = \mathbb{E}_q \mathbb{E}_{r \sim \pi_\theta(\cdot|q)} [R(q, r) \cdot w(q, r)], \quad (16)$$

where the weight factor is approximately log-linear in Evidence Gain:

$$\log w(q, r) \approx \Delta(q, r) + c, \quad (17)$$

for some constant  $c$  that depends on the average variance  $\mathbb{E}_{q,r}[\text{Var}_e[\Delta_e]]$ .

This result confirms that the policy model’s intrinsic ICL ability provides an effective quality signal. Traces with higher Evidence Gain receive proportionally higher weights in the reweighted objective, without requiring any external evaluator. In this way, In-Context RLVR leverages the model’s own capacity to distinguish reasoning quality, enabling implicit reward reweighting through a simple modification to the training procedure.