

Towards Better RL Training Data Utilization via Second-Order Rollout

Anonymous ACL submission

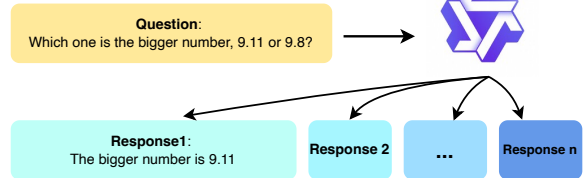
Abstract

Reinforcement Learning (RL) has empowered Large Language Models (LLMs) with strong reasoning capabilities, but vanilla RL mainly focuses on generation capability improvement by training with only first-order rollout (generating multiple responses for a question), and we argue that this approach fails to fully exploit the potential of training data because of the neglect of critique capability training. To tackle this problem, we further introduce the concept of second-order rollout (generating multiple critiques for a response) and propose a unified framework for jointly training generation and critique capabilities. Extensive experiments across various models and datasets demonstrate that our approach can utilize training data more effectively than vanilla RL and achieve better performance under the same training data. Additionally, we uncover several insightful findings regarding second-order rollout and critique training, such as the importance of label balance in critique training and the noise problem of outcome-based rewards, which can be mitigated through sampling techniques. Our work offers a preliminary exploration of dynamic data augmentation and joint generation-critique training in RL, providing meaningful inspiration for the further advancement of RL training.

1 Introduction

With the widespread application of Reinforcement Learning (RL) in post-training (Guo et al., 2025), Large Language Models (LLMs) have demonstrated remarkable reasoning capabilities, which inspires deeper investigations into RL training for LLMs. However, current RL training predominantly focuses on enhancing generation capability, often neglecting the development of critique capability, which can be a performance bottleneck for further improvement. Saunders et al. (2022) also categorizes model capabilities into *Generation* and

1. First-order Rollout



2. Second-order Rollout

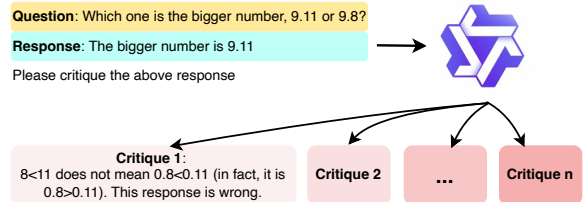


Figure 1: A demonstration of first/second-order rollout. The policy model generates multiple responses for a question in first-order rollout, and generates multiple critiques for a response in second-order rollout.

*Critique*¹: (1) *Generation* refers to the ability to produce a correct response to a given question; (2) *Critique* denotes the capacity to judge whether a response is correct and to identify specific errors in a wrong response. Intuitively, these two capabilities are not independent: Wang et al. (2025c) finds that fine-tuning a model using only critique data, even without any explicit generation data, can significantly improve its generation performance. Similarly, better critique performance with only generation training is also witnessed by Wang et al. (2025b). The neglect of critique training (Yu et al., 2025b; Xie et al., 2025) makes us wonder whether current RL training solely on generation capability can fully exploit the potential of training data, and we would like to further explore a joint RL training framework (Ruan et al., 2025; Wang et al., 2025a) for better data utilization.

¹The original classification delineates three distinct capabilities: Generation, Discrimination, and Critique. For simplicity, we subsume both Discrimination and Critique under the single term Critique.

In §3 we propose **Generation and Critique RL (GC-RL)** that jointly trains two capabilities with only generation training data by introducing the concept of *second-order rollout*. In vanilla RL, a policy model samples multiple responses for a given question during training, which we define as *first-order rollout*. Building on this, we further define *second-order rollout* as the process in which the policy model generates multiple critiques for a <question, response> pair, and these two processes are illustrated in Figure 1. At each training step, we sample questions from the training set for first-order rollout and <question, response> pairs from a data cache for second-order rollout, and these rollout is then combined to update the policy model collectively. Meanwhile, responses generated from first-order rollout are filtered and added to the data cache for subsequent training. It is worth noting that the second-order rollout can proceed naturally based on the results of the first-order rollout, and no additional training data is required, which can be viewed as a "free lunch" to some extent.

In §4, extensive experiments are conducted across different models and datasets, demonstrating that our GC-RL shows better data utilization and outperforms vanilla RL in both generation and critique capabilities under the same training data. Further we conduct more experiments to explore second-order rollout and critique training in §5 and find: 1. the data filter is crucial for maintaining balanced critique training (§5.1); 2. outcome-based reward is noisy for critique training and denoising can be achieved through multiple samplings (§5.2); 3. static data performs better in critique-only training and dynamic data are more suitable for joint training (§5.3); 4. fine-grained model critique behavior manipulation can be achieved through reward function adjustment. (§5.4).

Our contributions can be summarized as follows:

1. We introduce the concept of *second-order rollout* and propose GC-RL framework, which achieves better RL training data utilization by generation-critique joint training.
2. We conduct extensive experiments to show the effectiveness of our approach and draw some instructive conclusions about critique training.
3. Our work offers a preliminary exploration of dynamic data augmentation in RL training, providing meaningful inspiration for further advancement of RL.

2 Related Work

RL for LLMs The exceptional reasoning capabilities demonstrated by Deepseek-R1 (Guo et al., 2025) highlight the significant role of RL (Zhang et al., 2025) with verifiable reward in training LLMs. Beyond rule-based rewards, other forms such as model-based (Xu et al., 2025; Shao et al., 2025) and rubric-based (Gunjal et al., 2025; Huang et al., 2025) rewards can also be leveraged for RL training of LLMs. Another line of research focuses on developing RL algorithms suited for LLM training, including works like PPO (Schulman et al., 2017), GRPO (Shao et al., 2024), DAPO (Yu et al., 2025a), and GSPO (Zheng et al., 2025). There are also works exploring new RL training tasks and objectives: for instance, She et al. (2025) employs RL to train models in reconstructing questions from responses, while Dong et al. (2025) applies RL to enhance next-token prediction. Different from previous works, our work explores better RL training utilization through joint training of generation and critique capabilities.

LLM Critique The ability to provide critique constitutes a crucial component of LLM capabilities. High-quality critiques enable LLMs to perform self-correction (Pan et al., 2024; Yang et al., 2025a,b) more effectively, and can also enhance the reward signals produced by reward models (Yu et al., 2025c; Ankner et al., 2024; Ye et al., 2025) when incorporated into the context. Sun et al. (2024) proposes a framework for evaluating the quality of critiques, while other research efforts have focused on improving critique abilities through Supervised Fine-Tuning (Wang et al., 2025c), Direct Preference Optimization (Yu et al., 2025b) and RL (Xi et al., 2025; Xie et al., 2025; Tang et al., 2025). In contrast to prior work, our approach integrates critique training into the vanilla RL framework.

Data Augmentation Data augmentation (Wang et al., 2024; Chai et al., 2026) is an effective approach to enhance model performance in LLM training. Techniques such as back translation (Sennrich et al., 2016; Kulháněk et al., 2021) and rephrasing (Lu and Lam, 2023) can enrich the dataset without altering semantic meanings. Alternatively, another augmentation strategy involves leveraging LLMs to generate new data in zero-shot (Oh et al., 2023; Ubani et al., 2023) or in-context learning (Dai et al., 2025) settings. Unlike con-

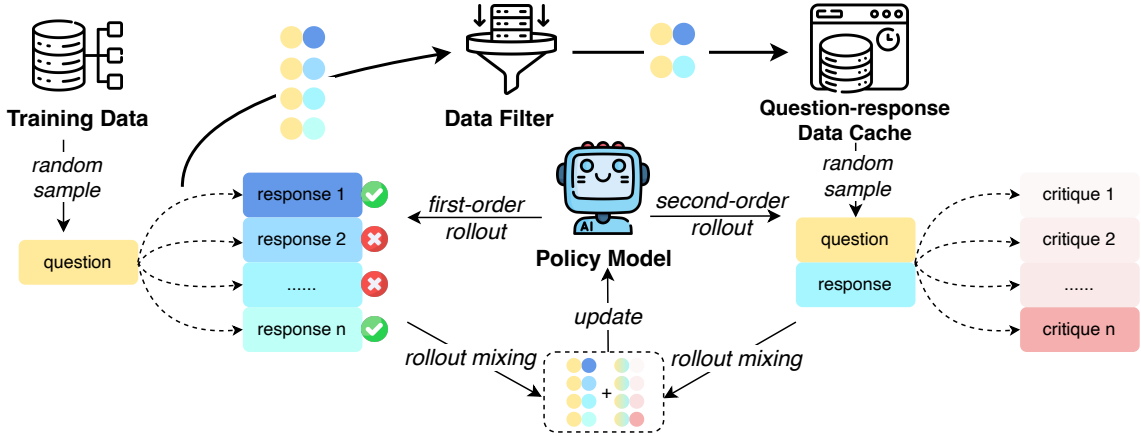


Figure 2: Flowchart of a single training step in GC-RL. First, a batch of questions is sampled from the training data, and multiple responses are generated (first-order rollout). Then, without replacement, a batch of \langle question, response \rangle pairs is sampled from the *Question-Response Data Cache*, and multiple critiques are generated (second-order rollout). These rollouts are combined and utilized jointly to update the policy model. In addition, the *Question-Response Data Cache* is maintained by processing the first-order rollouts through a *Data Filter* and adding the filtered data into the cache.

161 ventional offline data augmentation methods, our
 162 work can be viewed as an online data augmentation
 163 process conducted in RL.

164 3 Methodology

165 We introduce **Generation and Critique RL (GC-**
 166 **RL)**, an RL framework for jointly training the gener-
 167 ation and critique capabilities, and the overview
 168 of our approach is illustrated in the figure 2. At
 169 each RL training step, a batch of data is sampled
 170 from the initial training set, and the policy model
 171 generates multiple responses for each question (in
 172 conventional RL training, these responses would
 173 be directly used to update the policy model). We
 174 go further by feeding these \langle question, response \rangle
 175 pairs into a *data filter*, retaining a subset to be
 176 stored in a *question-response data cache*. Then
 177 a batch of \langle question, response \rangle pairs is sampled
 178 from the data cache, for which the policy model
 179 further generates critiques. The responses and cri-
 180 tiques obtained from these two rollout processes
 181 are combined, assigned rewards and advantages,
 182 and utilized to update the policy model.

183 **RL with Second-order Rollout** For each ques-
 184 tion q_i in the training set D , the policy model
 185 samples n responses r_1, r_2, \dots, r_n during vanilla
 186 RL training, a process we refer to as *first-order*
 187 *rollout*. For a specific \langle question, response \rangle pair
 188 $\langle q_i, r_i \rangle$, the policy model further samples n cri-
 189 tiques c_1, c_2, \dots, c_n for subsequent training, which
 190 we define as *second-order rollout*. In essence, the

191 first-order rollout aims to enhance the generation
 192 capability, enabling model to produce appropri-
 193 ate answers to given questions. The second-order
 194 rollout is designed to improve critique capability,
 195 empowering it to identify potential issues or errors
 196 within a response. It is worth noting that the second-
 197 order rollout can proceed naturally based on the
 198 results of the first-order rollout, and no additional
 199 training data is required, which can be viewed as a
 200 "free lunch" to some extent.

201 **Critique Data Filter** During the RL training
 202 process, the \langle question, response \rangle pairs obtained
 203 from the first-order rollout are passed through a
 204 *Data Filter*, which controls which pairs to retain.
 205 For a question q and its corresponding responses
 206 r_1, r_2, \dots, r_n , if all n responses are either entirely
 207 correct or entirely incorrect, all these data are dis-
 208 carded. Otherwise, one correct response r_{correct} and
 209 one incorrect response r_{wrong} will be randomly se-
 210 lected, and these two resulting pairs $\langle q, r_{\text{correct}} \rangle$
 211 and $\langle q, r_{\text{wrong}} \rangle$ will then stored in the data cache.
 212 Although simple, this data filter plays a crucial
 213 role in training stability, and removing this filter
 214 would introduce significant issues: (1). *imbalance*
 215 *between critique data and generation data*: for a
 216 single question q , the first-order rollout produces
 217 n responses and the second-order rollout further
 218 generates n^2 critiques, and critique data outnumber
 219 generation data by a factor of n ; (2). *imbalance*
 220 *within critique data*: we have observed that incor-
 221 rect responses tend to dominate over correct ones in

the first-order rollout, which can cause a label imbalance problem for further critique training, and further discussion is provided in §5.1.

Mixed Training During the update phase of the policy model, we perform training with a mixture of first-order rollout (responses) and second-order rollout (critiques). For each response r , a rule-based verifier is employed to check its correctness, with the reward function defined as:

$$R(r) = \begin{cases} 1, & r \text{ is correct} \\ 0, & r \text{ is wrong} \end{cases} \quad (1)$$

For each critique c generated based on $\langle \text{question}, \text{response} \rangle$ pair $\langle q, r \rangle$, the final judgment regarding the correctness of r is extracted, denoted as $Ext(c) \in \{correct, wrong\}$. Since intermediate steps of a critique are hard to verify (Sun et al., 2024), we only assign an outcome reward based on the final binary judgment, and the corresponding reward function is:

$$R(c) = \begin{cases} 0.7, & Ext(c) = correct \ \& \ R(r)=1 \\ 0.7, & Ext(c) = wrong \ \& \ R(r)=0 \\ 0, & else \end{cases} \quad (2)$$

For responses and critiques in the mixed rollout, the respective rewards are computed with the above reward functions, and corresponding advantages can then be obtained with GRPO (Shao et al., 2024) algorithm for further update of policy model.

4 Experiments

4.1 Experimental Setup

Models Our experiments are primarily conducted on the Qwen2.5 series (Yang et al., 2024), including Qwen2.5-(1.5B, 3B, 7B)-Base, to demonstrate the effectiveness of our method across models of varying scales. Additional experiments are performed on Llama-3.1-8B-Instruct (Grattafiori et al., 2024) and Mistral-7B-Instruct-v0.3 (Rastogi et al., 2025) to further validate the general applicability of our approach to different model architectures.

Dataset The training and evaluation mainly focus on mathematical reasoning tasks. DAPO-MATH-17k (Yu et al., 2025a) is employed as the primary training dataset: 1k data are randomly selected to construct cold-start data, while the remaining 16k data are utilized for RL training. Several established mathematical reasoning benchmarks are utilized for evaluation,

including Math-500 (Hendrycks et al., 2021), GSM8k (Cobbe et al., 2021), Minerva (Lewkowycz et al., 2022), AMC23, and OlympiadBench (He et al., 2024).

Implementation Details We employ the verl² framework for RL training and adopt GRPO (Shao et al., 2024) algorithm, with training hyper-parameters detailed in the Appendix A. All training and evaluation experiments are conducted on NVIDIA H20 96GB GPUs, and training is conducted on 8GPUs and evaluation is on a single GPU.

4.2 RL Training and Experimental Results

Cold Start Applying RL directly to models can lead to several issues: (1). Following the format of CFT (Wang et al., 2025c), we instruct the model to append a correctness judgment of its response at the end of the critique. However, the base models exhibit weak instruction-following capability, often failing to generate critiques that adhere to the required format. (2). Due to their limited reasoning performance, the intermediate reasoning steps within the generated critiques are often of low quality. To address these problems, we first distill 1,885 initial critique data from GPT-5³, and the prompt used for critique distillation is provided in the Appendix A. After filtering out data with incorrect formatting or erroneous final judgments, we obtain 1,339 high-quality training examples. Before starting RL, we perform Supervised Fine-Tuning (SFT) on this curated critique dataset to equip the model with a preliminary critique capability.

Baselines In addition to presenting the evaluation results of our **GC-RL** (Generation and Critique RL) approach, we also provide the outcomes of several baseline methods: (1) after cold start without RL training; (2) **G-RL** (Generation RL): vanilla RL training that performs only first-order rollout to enhance generation capability; (3) **C-RL** (Critique RL): for each question, 10 responses are sampled in advance to construct a balanced training set consisting of $\langle \text{question}, \text{response} \rangle$ pairs, and only second-order rollout is performed to train critique capability of LLMs. All these RL training are performed under the same training data.

Critique Evaluation In addition to generation capability evaluation, we also try to examine the cri-

²<https://github.com/volcengine/verl>

³gpt-5-chat-2025-08-07

Models	Methods	Math-500	GSM8k	Minerva	AMC23	Olympiad	Avg
		Generation Accuracy (%)					
Qwen2.5-7B	w/o RL	55.6	77.9	16.9	35.0	22.8	41.6
	C-RL	65.1	83.7	19.2	47.5	26.1	48.3
	G-RL	75.4	89.7	24.6	60.0	33.7	56.7
	GC-RL	77.6	92.0	24.6	62.5	39.8	59.3
Qwen2.5-3B	w/o RL	20.4	29.4	4.0	5.0	7.6	13.3
	C-RL	30.4	50.0	4.4	12.5	11.0	21.7
	G-RL	57.8	79.5	12.9	27.5	24.5	40.4
	GC-RL	61.8	81.4	14.9	32.5	25.2	43.2
Qwen2.5-1.5B	w/o RL	11.0	14.6	1.8	5.0	3.3	7.1
	C-RL	21.3	25.6	3.6	10.0	5.9	13.3
	G-RL	45.1	67.2	8.7	25.0	12.7	31.7
	GC-RL	47.2	69.0	10.3	27.5	15.3	33.9
		Critique Accuracy (%)					
Qwen2.5-7B	C-RL	80.5	82.5	62.3	71.4	72.5	73.8
	GC-RL	84.6	88.3	67.1	79.4	73.8	78.6
Qwen2.5-3B	C-RL	67.4	66.7	57.1	64.7	61.2	63.4
	GC-RL	70.2	72.8	57.1	66.2	61.5	65.6
Qwen2.5-1.5B	C-RL	59.5	57.7	53.8	58.8	55.0	57.0
	GC-RL	61.4	60.6	57.2	61.3	57.5	59.6

Table 1: Generation and critique capabilities evaluation results on Qwen-2.5-(1.5B,3B,7B). GC-RL outperforms all other RL training methods in both generation and critique capabilities.

312 critique capability after RL training. To construct the
313 evaluation dataset, we utilize 5 datasets in genera-
314 tion evaluation as seed data, and sample responses
315 with Qwen2.5-(1.5B,7B,72)-Instruct, respectively
316 (10 responses for each model). The final answer
317 is required to be enclosed in `boxed{}`, and we filter
318 out responses that dissatisfy this format require-
319 ment, as well as questions for which all sampled
320 responses are either entirely correct or entirely in-
321 correct. From the remaining data, we randomly
322 select one correct response and an incorrect one
323 for each question, discarding all other responses,
324 and this process yields critique evaluation datasets
325 in which the correct and incorrect responses are
326 1:1. Since assessing the accuracy of the interme-
327 diate reasoning steps within a critique is particularly
328 challenging, we focus solely on evaluating whether
329 the final judgment of critique on the correctness of
330 the response is accurate, which essentially reduces
331 the task to a binary classification.

332 **Results** The evaluation experiments are con-
333 ducted on both generation and critique capabilities:
334 For generation tasks, model-generated answers are
335 verified by comparing them with reference answers
336 and the final accuracy is reported; For critique tasks,
337 we extract the generated final judgment on the cor-
338 rectness of a response and measure the binary clas-
339 sification accuracy. We show experimental results

340 for Qwen-2.5-(1.5B,3B,7B) in Table 1, and pro-
341 vide more results on Llama-3.1-8B-instruct and
342 Mistral-7B-Instruct-v0.3 in Appendix B, finding
343 that: 1. Even if a model is trained solely with C-RL
344 (without generation training), its generation ability
345 significantly improves—though such enhancement
346 remains far inferior to that achieved through G-RL.
347 2. After training with GC-RL, the model attains op-
348 timal performance in both generation and critique
349 capabilities. On one hand, its generation capability
350 surpasses that of models trained with G-RL; on the
351 other hand, its critique ability exceeds that of mod-
352 els trained with C-RL. This result suggests a certain
353 coupling between critique and generation capabil-
354 ities, and further demonstrates that joint training
355 yields superior overall performance compared to
356 training for each capability independently.

357 5 Analysis

358 First-order rollout in RL training has been thor-
359 oughly studied by previous works, so we conduct a
360 more detailed analysis of second-order rollout and
361 critique training in this section. First, we discuss
362 the label imbalance problem in critique training and
363 demonstrate the effectiveness of our data filter both
364 theoretically and experimentally (§5.1). Next, we
365 discuss the reward noise problem in critique train-
366 ing and explore a sampling-based denoising strat-

		Math-500	GSM8k	Minerva	AMC23	Olympiad	Avg
Generation	Random Sampling	75.0	90.9	22.9	60.0	37.5	57.3
	Random Sampling + Reweight	77.2	91.5	23.2	60.0	38.2	58.0
	Data Filter	77.6	92.0	24.6	62.5	39.8	59.3
Critique	Random Sampling	82.3	84.0	63.3	73.5	71.6	74.9
	Random Sampling + Reweight	83.7	86.7	65.8	77.9	72.8	77.4
	Data Filter	84.6	88.3	67.1	79.4	73.8	78.6

Table 2: A comparison of model performance on Qwen2.5-7B when sampling responses with/without a data filter. Random sampling leads to the worst performance, though adding reward reweighting can alleviate this issue. Utilizing the data filter achieves the best performance.

egy (§5.2). We then compare static and dynamic data in critique training, observing that dynamic data is more suitable for GC-RL, while static data works better for C-RL (§5.3). Finally, by adjusting the reward function, we achieve fine-grained critique behavior manipulation of LLMs after RL training (§5.4).

5.1 Towards Balanced Critique Training

We theoretically analyze why GC-RL without a data filter can lead to training data imbalance and restrict the critique capability of LLMs. To mitigate this problem, we explore employing a reward reweighting strategy and utilizing a data filter, and finally conduct comparative experiments to validate the effectiveness of our data filter.

Alleviating imbalance problem with reward reweighting The ultimate judgement of whether a response is correct or not in a critique is essentially a binary classification task, but the number of erroneous responses significantly outweighs the correct ones in the first-order rollout. Intuitively, this label imbalance problem can impact subsequent critique training, and we also theoretically analyze the effect of data imbalance on the performance of the critique training, with detailed discussions shown in Appendix C. To mitigate this problem, we also explore reweighting and scaling rewards for positive and negative data to balance their contributions, and the weighted reward function is defined as:

$$R_w(c) = \begin{cases} \frac{0.7}{E[R(r)]}, & Ext(c) = correct \ \& \ R(r)=1 \\ \frac{0.7}{1-E[R(r)]}, & Ext(c) = wrong \ \& \ R(r)=0 \\ 0, & else \end{cases} \quad (3)$$

where $E[R(r)]$ is the expected reward for response r during RL training. It can be mathematically proved that the weighted reward is unbiased and

does not incentivize the model to judge an uncertain response as correct or wrong, and the detailed proof is shown in Appendix C. Essentially, this weighting strategy amplifies the reward for rare classes, enabling the model to learn more effectively from such data and thereby approximating balanced training.

Empirical comparison of training with/without a data filter

Comparison experiments are conducted on Qwen2.5-7B with GC-RL training under three settings: (1) randomly sampling responses without data filter, (2) randomly sampling responses and training with the weighted reward function in Equation 3, and (3) utilizing the data filter in §3. As the experimental results shown in Table 2, random sampling strategy leads to the worst generation and critique capabilities because of label imbalance, and this problem can be alleviated to some extent with a weighted reward function. Although reward weighting can achieve balance in the reward level for imbalanced data, utilizing a data filter can achieve inherently balanced training data and yield the best performance.

5.2 Reward Noise & Denoising in Critique RL

Reward noise problem in critique RL Obtaining precise rewards for critiques is more challenging than for responses. For instance, in mathematical problems, since we have a corresponding answer to each question, the correctness of a response can be easily rule-based verified and a precise reward can then be assigned accordingly. However, for critiques, it is difficult to verify the correctness of each intermediate step, and we can only assign rewards based on whether the final binary classification result is correct. Generating responses is essentially a generation task, and it is rare for a model to produce intermediate errors yet still arrive at the correct final answer. In contrast, generating critiques is essentially a binary classification

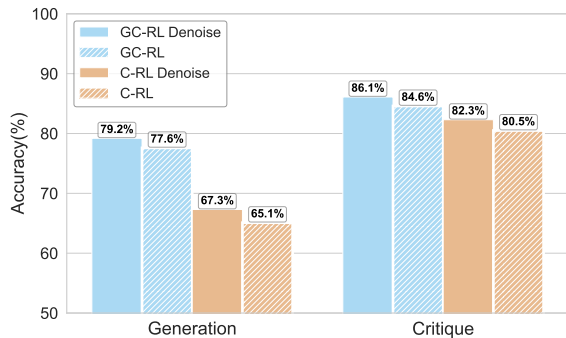


Figure 3: A comparison of model performance of Qwen2.5-7B with/without reward denoising strategy on Math-500. In both GC-RL and C-RL settings, reward denoising can improve model performance on both generation and critique capabilities.

task, and even random guessing can yield correct answers with a 50% probability, leading to many critiques with incorrect intermediate steps but correct final judgement. Ideally, for critiques with correct outcomes, we should differentiate between those with erroneous intermediate steps and those with correct ones, assigning lower and higher rewards, respectively. However, verifying intermediate steps is challenging in practice, so critiques with correct results are often given the same reward, and we refer to such rewards as *noised rewards*. Guo et al. (2025) has demonstrated significant success in Reinforcement Learning with Verifiable Rewards (RLVR), which fundamentally relies on *oracle rewards*. Liu et al. (2025); Shi and Jin (2025); Whitehouse et al. (2025) also attempt RL on classification tasks, showing that even with *noised rewards*, model performance can be improved to some extent, which is also witnessed by our experiments in §4.

Exploring critique reward denoising based on self-correction with multiple samplings In LLM self-correction process (Kamoi et al., 2024; Pan et al., 2024), a model is provided with three components: <question, response, critique>, and then instructed to modify the original response based on the potential issues identified by the critique and generate a refined response. Intuitively, the higher the quality of the critique—i.e., the more accurately it identifies problems in the original response—the greater the likelihood that the refined response will be correct. Thus, we can inversely estimate the quality of a critique based on the quality of its corresponding refined response. Similar to (Tang et al., 2025; Xie et al., 2025; Yu et al.,

2025b), for a given critique, we allow the model to perform self-correction and sample n refined responses, with the number of correct ones denoted as k , based on which we then propose the following reward function to estimate critique quality:

$$R_q(c) = \begin{cases} 0.1 * \frac{k}{n}, & Ext(c) == correct \\ 0.7 * \frac{k}{n}, & Ext(c) == wrong \end{cases} \quad (4)$$

Although it is difficult to directly verify the correctness of intermediate steps in critique, this sampling method can indirectly estimate it to some extent, thereby reducing noise in the reward. We utilize the sum of the outcome reward from Equation 2 and the estimated reward obtained from Equation 4, denoted as $R(c) + R_q(c)$, as the final reward for critique, and compare it with using only the outcome reward $R(c)$. Theoretically, a larger number of samples n leads to better noise reduction and more accurate reward values, but at a higher computational cost. To make the computational overhead controllable, experiments are conducted with $n = 1$, and the experimental results on Math-500 are shown in Figure 3 (with more results shown in Figure 9 in Appendix B). A performance improvement in both generation and critique capabilities can be witnessed under both CG-RL and C-RL settings with our reward denoising strategy.

5.3 Static vs. Dynamic Data

A comparison of static and dynamic responses during critique RL training. In our approach, dynamic self-generated responses are utilized when generating second-order rollout, and an alternative strategy involves utilizing static responses (Ruan et al., 2025; Xie et al., 2025) for critique RL training, where pre-prepared <question, response> pairs remain fixed throughout the RL process. A performance comparison of static and dynamic training data is conducted under GC-RL and C-RL settings for both generation and critique capabilities on Qwen2.5-7B, and the experimental results on Math-500 are presented in Figure 4 (with more results shown in Figure 7 in Appendix B). We find that training with dynamically self-generated response data leads to higher performance in both generation and critique capabilities under GC-RL setting. However, under C-RL setting, where responses are abandoned and only critiques are utilized to update the policy model, dynamic data suffers from a severe reward hacking problem, and the static data strategy significantly outperforms the dynamic

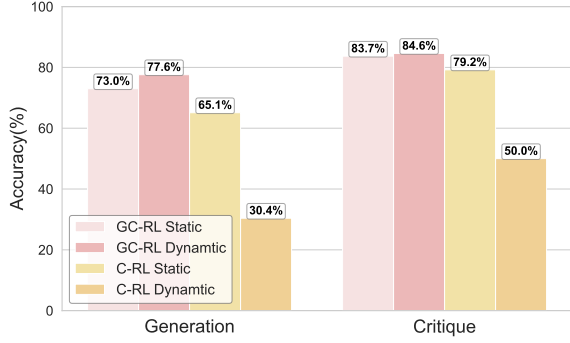


Figure 4: Performance of Qwen2.5-7B on Math-500 under GC-RL and C-RL settings with static and dynamic critique training data. Dynamic data outperforms static data in the GC-RL setting, while the opposite holds for C-RL.

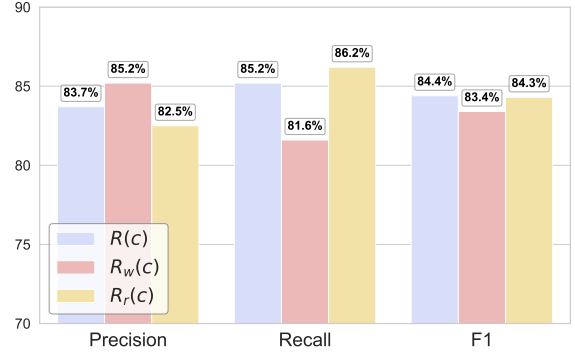


Figure 5: A comparison of critique performance of Qwen2.5-7B with different reward function on Math-500. Compared to baseline $R(c)$, $R_w(c)$ leads to a higher precision while $R_r(c)$ generates a higher recall.

523 data. With dynamic data, the model seems to identify
 524 a shortcut to maximize rewards during RL: it
 525 deliberately generates incorrect responses in the
 526 generation stage (though producing a correct re-
 527 sponse is challenging, generating an incorrect one
 528 can be quite easy), then it labels all responses as
 529 incorrect to get the reward in the critique stage. To
 530 summarize, dynamic data is more suitable for GC-
 531 RL training, while static data is more appropriate
 532 for C-RL training.

5.4 Fine-grained Critique Behavior Manipulation

533 There are often different requirements for the bi-
 534 nary classification performance in different scenar-
 535 ios. For instance, a higher recall rate is demanded
 536 in disease screening, while a higher precision rate
 537 is required in recommendation systems. In the reward
 538 function defined in Equation 2, the same reward is
 539 assigned to both correct and incorrect responses as
 540 long as the critique identifies them correctly, and
 541 we also explore training with more fine-grained
 542 reward functions to better control the critique be-
 543 havior of LLMs. For example, to encourage the
 544 model to be more inclined to classify a response as
 545 incorrect when it is uncertain, we try the following
 546 reward function:
 547
 548

$$R_w(c) = \begin{cases} 0.6, & Ext(c) = correct \ \& \ R(r)=1 \\ 0.8, & Ext(c) = wrong \ \& \ R(r)=0 \\ 0, & else \end{cases} \quad (5)$$

549 Conversely, to steer the model toward classifying
 550 uncertain responses as correct, we assign a large
 551 reward value to correct responses and employ the
 552 following reward function:
 553

$$R_r(c) = \begin{cases} 0.8, & Ext(c) = correct \ \& \ R(r)=1 \\ 0.6, & Ext(c) = wrong \ \& \ R(r)=0 \\ 0, & else \end{cases} \quad (6)$$

554
 555 The experimental results on Math-500 of uti-
 556 lizing $R_w(c)$ and $R_r(c)$ are presented in Figure
 557 5, and more results can be found in Figure 8 in
 558 Appendix B). Compared with the baseline $R(c)$,
 559 we observe that when $R_w(c)$ is applied, the model
 560 achieves higher precision but lower recall; whereas
 561 with $R_r(c)$, the precision decreases while recall
 562 increases. By adjusting the reward function, we
 563 can exert finer-grained control over the critique
 564 behavior of LLMs.

6 Conclusion

565
 566 Based on the first-order rollout in vanilla RL (gen-
 567 erating multiple responses for a question), we fur-
 568 ther introduce the concept of second-order rollout
 569 (generating multiple critiques for a response) and
 570 propose GC-RL, a unified framework to train gen-
 571 eration and critique capabilities jointly. Extensive
 572 experiments across various models and datasets
 573 demonstrate that our approach can more effectively
 574 utilize training data compared to vanilla RL, achiev-
 575 ing superior performance under the same training
 576 data. Additionally, we uncover some insightful
 577 findings related to second-order rollout and critique
 578 training, such as the importance of label balance in
 579 critique training. Our work serves as an initial ex-
 580 ploration into dynamic data augmentation and joint
 581 training of generation and critique in RL training,
 582 offering meaningful insights for further advance-
 583 ments in RL for LLMs.

584 Limitations

585 Our work represents a preliminary attempt to inte-
586 grate dynamic data augmentation with joint train-
587 ing of generation and critique into RL, and still ex-
588 hibits several limitations that warrant further explo-
589 ration. For simplicity, we only employ the GRPO
590 algorithm for RL training, and the applicability
591 of our approach to other RL algorithms (such as
592 PPO) remains to be investigated. Moreover, exper-
593 iments are confined to the mathematical domain
594 using models with fewer than 10B parameters, and
595 extending our approach to larger-scale RL training
596 with multi-domain training data and larger models
597 is considered to be an important direction for future
598 work. Additionally, we have observed that our ap-
599 proach converges more slowly compared to vanilla
600 RL, essentially trading computational resources
601 for improved performance. Our approach also re-
602 quires that responses can be rule-based verified,
603 making it less straightforward to apply to RL tasks
604 where responses are harder to evaluate (e.g., rubric-
605 based RL training (Gunjal et al., 2025; Huang et al.,
606 2025)). How to perform second-order rollout and
607 assign rewards to critiques in such settings is also
608 worth further exploration.

609 Ethical Considerations

610 The data utilized are open for research, and LLMs
611 in the experiments are all publicly available by
612 either parameters or API calls. Therefore, we do
613 not anticipate any ethical concerns in our research.

614 References

615 Zachary Ankner, Mansheej Paul, Brandon Cui,
616 Jonathan D Chang, and Prithviraj Ammanabrolu.
617 2024. Critique-out-loud reward models. *arXiv*
618 *preprint arXiv:2408.11791*.

619 Yaping Chai, Haoran Xie, and Joe S Qin. 2026. Text
620 data augmentation for large language models: A com-
621 prehensive survey of methods, challenges, and oppor-
622 tunities. *Artificial Intelligence Review*, 59(1):35.

623 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
624 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
625 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
626 Nakano, and 1 others. 2021. Training verifiers
627 to solve math word problems. *arXiv preprint*
628 *arXiv:2110.14168*.

629 Haixing Dai, Zhengliang Liu, Wenxiong Liao, Xiaoke
630 Huang, Yihan Cao, Zihao Wu, Lin Zhao, Shaochen
631 Xu, Fang Zeng, Wei Liu, and 1 others. 2025. Auggpt:
632 Leveraging chatgpt for text data augmentation. *IEEE*
633 *Transactions on Big Data*.

Qingxiu Dong, Li Dong, Yao Tang, Tianzhu Ye, Yutao
Sun, Zhifang Sui, and Furu Wei. 2025. Reinforce-
ment pre-training. *arXiv preprint arXiv:2506.08007*. 634
635
636

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,
Abhinav Pandey, Abhishek Kadian, Ahmad Al-
Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,
Alex Vaughan, and 1 others. 2024. The llama 3 herd
of models. *arXiv preprint arXiv:2407.21783*. 637
638
639
640

Anisha Gunjal, Anthony Wang, Elaine Lau, Vaskar
Nath, Yunzhong He, Bing Liu, and Sean M. Hendryx.
2025. Rubrics as rewards: Reinforcement learning
beyond verifiable domains. In *NeurIPS 2025 Work-*
shop on Efficient Reasoning. 642
643
644
645
646

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao
Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi-
rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.
Deepseek-r1: Incentivizing reasoning capability in
llms via reinforcement learning. *arXiv preprint*
arXiv:2501.12948. 647
648
649
650
651
652

Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu,
Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie
Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan
Liu, and Maosong Sun. 2024. OlympiadBench:
A challenging benchmark for promoting AGI with
olympiad-level bilingual multimodal scientific prob-
lems. In *Proceedings of the 62nd Annual Meeting of*
the Association for Computational Linguistics (Vol-
ume 1: Long Papers), pages 3828–3850, Bangkok,
Thailand. Association for Computational Linguistics. 653
654
655
656
657
658
659
660
661
662

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul
Arora, Steven Basart, Eric Tang, Dawn Song, and
Jacob Steinhardt. 2021. Measuring mathematical
problem solving with the MATH dataset. In *Thirty-*
fifth Conference on Neural Information Processing
Systems Datasets and Benchmarks Track (Round 2). 663
664
665
666
667
668

Zenan Huang, Yihong Zhuang, Guoshan Lu, Zeyu Qin,
Haokai Xu, Tianyu Zhao, Ru Peng, Jiaqi Hu, Zhan-
ming Shen, Xiaomeng Hu, and 1 others. 2025. Re-
inforcement learning with rubric anchors. *arXiv*
preprint arXiv:2508.12790. 669
670
671
672
673

Ryo Kamoi, Yusen Zhang, Nan Zhang, Jiawei Han,
and Rui Zhang. 2024. When can LLMs actually
correct their own mistakes? a critical survey of self-
correction of LLMs. *Transactions of the Association*
for Computational Linguistics, 12:1417–1440. 674
675
676
677
678

Jonáš Kulhánek, Vojtěch Hudeček, Tomáš Nekvinda,
and Ondřej Dušek. 2021. AuGPT: Auxiliary tasks
and data augmentation for end-to-end dialogue with
pre-trained language models. In *Proceedings of the*
3rd Workshop on Natural Language Processing for
Conversational AI, pages 198–210, Online. Associa-
tion for Computational Linguistics. 679
680
681
682
683
684
685

Aitor Lewkowycz, Anders Johan Andreassen,
David Dohan, Ethan Dyer, Henryk Michalewski,
Vinay Venkatesh Ramasesh, Ambrose Slone, Cem
Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu,
Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 686
687
688
689
690

691	2022. Solving quantitative reasoning problems with language models . In <i>Advances in Neural Information Processing Systems</i> .	747
692		748
693		749
694	Zijun Liu, Peiyi Wang, Runxin Xu, Shirong Ma,	750
695	Chong Ruan, Peng Li, Yang Liu, and Yu Wu. 2025.	751
696	Inference-time scaling for generalist reward model-	
697	ing, 2025. URL https://arxiv.org/abs/2504.02495 .	
698	Hongyuan Lu and Wai Lam. 2023. Epa: easy prompt	
699	augmentation on large language models via multi-	
700	ple sources and multiple targets. <i>arXiv preprint</i>	
701	<i>arXiv:2309.04725</i> .	
702	Seokjin Oh, Su Ah Lee, and Woohwan Jung. 2023.	
703	Data augmentation for neural machine translation	
704	using generative language model. <i>arXiv preprint</i>	
705	<i>arXiv:2307.16833</i> .	
706	Liangming Pan, Michael Saxon, Wenda Xu, Deepak	
707	Nathani, Xinyi Wang, and William Yang Wang. 2024.	
708	Automatically correcting large language models: Sur-	
709	veying the landscape of diverse automated correction	
710	strategies . <i>Transactions of the Association for Com-</i>	
711	<i>putational Linguistics</i> , 12:484–506.	
712	Abhinav Rastogi, Albert Q Jiang, Andy Lo, Gabrielle	
713	Berrada, Guillaume Lample, Jason Rute, Joep Bar-	
714	mentlo, Karmesh Yadav, Kartik Khandelwal, Khy-	
715	athi Raghavi Chandu, and 1 others. 2025. <i>Magistral</i> .	
716	<i>arXiv preprint arXiv:2506.10910</i> .	
717	Chi Ruan, Dongfu Jiang, Yubo Wang, and Wenhua	
718	Chen. 2025. Critique-coder: Enhancing coder models	
719	by critique reinforcement learning. <i>arXiv preprint</i>	
720	<i>arXiv:2509.22824</i> .	
721	William Saunders, Catherine Yeh, Jeff Wu, Steven Bills,	
722	Long Ouyang, Jonathan Ward, and Jan Leike. 2022.	
723	Self-critiquing models for assisting human evaluators .	
724	<i>CoRR</i> , abs/2206.05802.	
725	John Schulman, Filip Wolski, Prafulla Dhariwal,	
726	Alec Radford, and Oleg Klimov. 2017. Proxi-	
727	mal policy optimization algorithms. <i>arXiv preprint</i>	
728	<i>arXiv:1707.06347</i> .	
729	Rico Sennrich, Barry Haddow, and Alexandra Birch.	
730	2016. Improving neural machine translation models	
731	with monolingual data . In <i>Proceedings of the 54th</i>	
732	<i>Annual Meeting of the Association for Computational</i>	
733	<i>Linguistics (Volume 1: Long Papers)</i> , pages 86–96,	
734	Berlin, Germany. Association for Computational Lin-	
735	guistics.	
736	Zhihong Shao, Yuxiang Luo, Chengda Lu, ZZ Ren,	
737	Jiewen Hu, Tian Ye, Zhibin Gou, Shirong Ma, and	
738	Xiaokang Zhang. 2025. Deepseekmath-v2: To-	
739	wards self-verifiable mathematical reasoning. <i>arXiv</i>	
740	<i>preprint arXiv:2511.22570</i> .	
741	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	
742	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	
743	Zhang, YK Li, Yang Wu, and 1 others. 2024.	
744	Deepseekmath: Pushing the limits of mathematical	
745	reasoning in open language models. <i>arXiv preprint</i>	
746	<i>arXiv:2402.03300</i> .	
	Shuaijie She, Yu Bao, Yu Lu, Lu Xu, Tao Li, Wenhao	747
	Zhu, Shujian Huang, Shanbo Cheng, Lu Lu, and Yux-	748
	uan Wang. 2025. Dupo: Enabling reliable llm self-	749
	verification via dual preference optimization. <i>arXiv</i>	750
	<i>preprint arXiv:2508.14460</i> .	751
	Wenlei Shi and Xing Jin. 2025. Heimdall: test-time	752
	scaling on the generative verification. <i>arXiv preprint</i>	753
	<i>arXiv:2504.10337</i> .	754
	Shichao Sun, Junlong Li, Weizhe Yuan, Ruifeng Yuan,	755
	Wenjie Li, and Pengfei Liu. 2024. The critique of	756
	critique . In <i>Findings of the Association for Computa-</i>	757
	<i>tional Linguistics: ACL 2024</i> , pages 9077–9096,	758
	Bangkok, Thailand. Association for Computational	759
	Linguistics.	760
	Qiaoyu Tang, Hao Xiang, Le Yu, Bowen Yu, Hongyu	761
	Lin, Yaojie Lu, Xianpei Han, Le Sun, and Jun-	762
	yang Lin. 2025. Refcritic: Training long chain-	763
	of-thought critic models with refinement feedback.	764
	<i>arXiv preprint arXiv:2507.15024</i> .	765
	Solomon Ubani, Suleyman Olcay Polat, and Rodney	766
	Nielsen. 2023. Zeroshotdataaug: Generating and aug-	767
	menting training data with chatgpt. <i>arXiv preprint</i>	768
	<i>arXiv:2304.14334</i> .	769
	Ke Wang, Jiahui Zhu, Minjie Ren, Zeming Liu, Shiwei	770
	Li, Zongye Zhang, Chenkai Zhang, Xiaoyu Wu, Qiqi	771
	Zhan, Qingjie Liu, and 1 others. 2024. A survey on	772
	data synthesis and augmentation for large language	773
	models. <i>arXiv preprint arXiv:2410.12896</i> .	774
	Xiaoxuan Wang, Bo Liu, Song Jiang, Jingzhou Liu,	775
	Jingyuan Qi, Xia Chen, and Baosheng He. 2025a.	776
	From solving to verifying: A unified objective	777
	for robust reasoning in llms. <i>arXiv preprint</i>	778
	<i>arXiv:2511.15137</i> .	779
	Xiyao Wang, Chunyuan Li, Jianwei Yang, Kai Zhang,	780
	Bo Liu, Tianyi Xiong, and Furong Huang. 2025b.	781
	Llava-critic-r1: Your critic model is secretly a strong	782
	policy model. <i>arXiv preprint arXiv:2509.00676</i> .	783
	Yubo Wang, Xiang Yue, and Wenhua Chen. 2025c. Cri-	784
	tique fine-tuning: Learning to critique is more effec-	785
	tive than learning to imitate . In <i>Second Conference</i>	786
	<i>on Language Modeling</i> .	787
	Chenxi Whitehouse, Tianlu Wang, Ping Yu, Xian	788
	Li, Jason Weston, Ilia Kulikov, and Swarnadeep	789
	Saha. 2025. J1: Incentivizing thinking in llm-as-	790
	a-judge via reinforcement learning. <i>arXiv preprint</i>	791
	<i>arXiv:2505.10320</i> .	792
	Zhiheng Xi, Jixuan Huang, Xin Guo, Boyang Hong,	793
	Dingwen Yang, Xiaoran Fan, Shuo Li, Zehui Chen,	794
	Junjie Ye, Siyu Yuan, and 1 others. 2025. Critique-	795
	r1: Training language models for critiquing through	796
	two-stage reinforcement learning. <i>arXiv preprint</i>	797
	<i>arXiv:2510.24320</i> .	798
	Zhihui Xie, Jie chen, Liyu Chen, Weichao Mao, Jingjing	799
	Xu, and Lingpeng Kong. 2025. Teaching language	800
	models to critique via reinforcement learning . In	801

802 *ICLR 2025 Third Workshop on Deep Learning for*
803 *Code.*

804 Zhangchen Xu, Yuetai Li, Fengqing Jiang, Bhaskar Ra-
805 masubramanian, Luyao Niu, Bill Yuchen Lin, and
806 Radha Poovendran. 2025. Tinyv: Reducing false neg-
807 atives in verification improves rl for llm reasoning.
808 *arXiv preprint arXiv:2505.14625.*

809 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng,
810 Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan
811 Li, Dayiheng Liu, Fei Huang, and 1 others.
812 2024. Qwen2 technical report. *arXiv preprint*
813 *arXiv:2407.10671.*

814 Zhe Yang, Yichang Zhang, Yudong Wang, Ziyao Xu,
815 Junyang Lin, and Zhifang Sui. 2025a. **Confidence v.s.**
816 **critique: A decomposition of self-correction capabil-**
817 **ity for LLMs.** In *Proceedings of the 63rd Annual*
818 *Meeting of the Association for Computational Lin-*
819 *guistics (Volume 1: Long Papers)*, pages 3998–4014,
820 Vienna, Austria. Association for Computational Lin-
821 guistics.

822 Zhe Yang, Yichang Zhang, Yudong Wang, Ziyao Xu,
823 Junyang Lin, and Zhifang Sui. 2025b. **A probabilis-**
824 **tic inference scaling theory for LLM self-correction.**
825 In *Proceedings of the 2025 Conference on Empirical*
826 *Methods in Natural Language Processing*, pages
827 13584–13598, Suzhou, China. Association for Com-
828 putational Linguistics.

829 Zihuiwen Ye, Fraser David Greenlee, Max Bartolo, Phil
830 Blunsom, Jon Ander Campos, and Matthias Gallé.
831 2025. **Improving reward models with synthetic cri-**
832 **tiques.** In *Findings of the Association for Computa-*
833 *tional Linguistics: NAACL 2025*, pages 4506–4520,
834 Albuquerque, New Mexico. Association for Compu-
835 tational Linguistics.

836 Qiyong Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xi-
837 aochen Zuo, YuYue, Weinan Dai, Tiantian Fan, Gao-
838 hong Liu, Juncai Liu, LingJun Liu, Xin Liu, Haibin
839 Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yux-
840 uan Tong, Chi Zhang, Mofan Zhang, and 17 others.
841 2025a. **DAPO: An open-source LLM reinforcement**
842 **learning system at scale.** In *The Thirty-ninth An-*
843 *ual Conference on Neural Information Processing*
844 *Systems.*

845 Tianshu Yu, Chao Xiang, Mingchuan Yang, Pei Ke, Bosi
846 Wen, Cunxiang Wang, Jiale Cheng, Li Zhang, Xinyu
847 Mu, Chuxiong Sun, and Minlie Huang. 2025b. **Train-**
848 **ing language model to critique for better refinement.**
849 In *Findings of the Association for Computational*
850 *Linguistics: ACL 2025*, pages 26760–26804, Vienna,
851 Austria. Association for Computational Linguistics.

852 Yue Yu, Zhengxing Chen, Aston Zhang, Liang Tan,
853 Chenguang Zhu, Richard Yuanzhe Pang, Yundi Qian,
854 Xuwei Wang, Suchin Gururangan, Chao Zhang,
855 Melanie Kambadur, Dhruv Mahajan, and Rui Hou.
856 2025c. **Self-generated critiques boost reward model-**
857 **ing for language models.** In *Proceedings of the 2025*
858 *Conference of the Nations of the Americas Chap-*
859 *ter of the Association for Computational Linguistics:*

Human Language Technologies (Volume 1: Long Pa-
860 *pers)*, pages 11499–11514, Albuquerque, New Mex-
861 ico. Association for Computational Linguistics. 862

Kaiyan Zhang, Yuxin Zuo, Bingxiang He, Youbang Sun,
863 Runze Liu, Che Jiang, Yuchen Fan, Kai Tian, Guoli
864 Jia, Pengfei Li, and 1 others. 2025. A survey of
865 reinforcement learning for large reasoning models.
866 *arXiv preprint arXiv:2509.08827.* 867

Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui
868 Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong
869 Liu, Rui Men, An Yang, and 1 others. 2025.
870 Group sequence policy optimization. *arXiv preprint*
871 *arXiv:2507.18071.* 872

Appendix 873

A More Implementation Details 874

We provide more implementation details in this
875 section. We show the prompt utilized for cold start
876 data distillation in Figure 6. We show some hyper-
877 parameters in RL training in Table 3. 878

Prompt for Critique Data Distillation

```
#Question#:
<insert question>
#Solution#:
<insert question>
#Instruction#:
Please verify step by step and judge whether the solu-
tion is correct, and end your answer with **Conclu-
sion: right/wrong [END]**
```

Figure 6: Our prompt fed to GPT-5 for critique data generation, which is comprised of a question, a solution, and critique instruction.

train batch size		512
ppo mini batch size		128
rollout n		5
adv estimator		grp
kl loss coef		1e-3
learning rate		1e-6
max prompt length		4096
max response length		4096
clip ratio		0.2
epochs		10

Table 3: RL training hyper-parameters.

B More Experimental Results 879

We show more experimental results in this section. 880
Main experiments on Mistral-7B-Instruct-v0.3 and 881

Llama-3.1-8B-Instruct are shown in Table 4. Performance comparison of static and dynamic critique training data is shown in Figure 7; performance comparison of different reward functions is shown in Figure 8; performance comparison of training with/without reward denoising strategy is shown in Figure 9.

C Discussions on Label Imbalance Problem

A theoretical analysis of data imbalance problem in critique training without the data filter. Theoretically, generating a critique can be viewed as a binary classification task, requiring the critique to include a final judgment on whether the response is correct or not, and binary classification tasks are susceptible to label distribution imbalance in the training data: when one category predominates, the trained model tends to favor predicting that category. In §3, the data processed through the first-order rollout and subsequently filtered by a data filter are stored in the *Question-response Data Cache* to support subsequent critique RL training, and this data filter ensures a 1:1 ratio between correct and incorrect responses, thereby preventing label bias in subsequent critique training. If the data filter is omitted and responses are sampled uniformly at random, in the first-order rollout, $E[R(r)] \times 100\%$ of responses are correct and $(1 - E[R(r)]) \times 100\%$ are incorrect, where E denotes the mathematical expectation. Consequently, the ratio of correct to incorrect responses under random sampling is also $E[R(r)] : (1 - E[R(r)])$. Similar to Yang et al. (2025a), let P_1 and P_2 denote the probabilities that the model correctly identifies correct and wrong responses, respectively, during the critique phase. On a validation set with a balanced 1:1 positive-to-negative ratio, the expected validation reward is $E[R_{val}(c)] = \frac{0.7P_1 + 0.7P_2}{2}$, and the expected reward for critiques during RL training is:

$$\begin{aligned}
& E[R(c)] \\
&= E[R(r)] * P_1 * 0.7 + (1 - E[R(r)]) * P_2 * 0.7 \\
&= 0.7E[R(r)]P_1 \\
&\quad + (1 - E[R(r)])(2E[R_{val}(c)] - 0.7P_1) \\
&= 0.7(2E[R(r)] - 1)P_1 \\
&\quad + 2(1 - E[R(r)])E[R_{val}(c)]
\end{aligned} \tag{7}$$

Though $E[R(r)]$ and $E[R_{val}(c)]$ gradually improve throughout the whole RL training process, they can be approximately viewed as constants be-

tween adjacent RL steps. From this point of view and Equation 7, we can see that the expected critique reward is in proportion to P_1 . Similar to Yang et al. (2025a), P_1 and P_2 also exhibit a competitive trade-off and $P_1 + P_2$ can be viewed as a constant between adjacent RL steps. When $2E[R(r)] - 1 > 0$, this reward encourages the model to increase P_1 and decrease P_2 ; when $2E[R(r)] - 1 < 0$, a decrease in P_1 raises the expected reward, thereby incentivizing the model to reduce P_1 and increase P_2 . For example, when training Qwen2.5-7B with GC-RL in §4, we find $E[R(r)] \approx 0.1$ in early training and $E[R(r)] \approx 0.45$ in late stage, consistently satisfying $2E[R(r)] - 1 < 0$ and incentivizing the model to decrease P_1 . Consequently, a model trained with randomly sampled responses for second-order rollout exhibits lower P_1 and higher P_2 , meaning it tends to classify responses as incorrect during critique.

The weighted reward function is unbiased We give a proof of the unbiasedness of the following weighted reward function:

$$R_w(c) = \begin{cases} \frac{0.7}{E[R(r)]}, & Ext(c) = correct \ \& \ R(r)=1 \\ \frac{0.7}{1-E[R(r)]}, & Ext(c) = wrong \ \& \ R(r)=0 \\ 0, & else \end{cases} \tag{8}$$

Under this scheme, the expected reward for critiques during RL training becomes:

$$\begin{aligned}
& E[R(c)] \\
&= E[R(r)]P_1 \frac{0.7}{E[R(r)]} \\
&\quad + (1 - E[R(r)])P_2 \frac{0.7}{1 - E[R(r)]} \\
&= 0.7(P_1 + P_2) \\
&= E[R_{val}(c)]
\end{aligned} \tag{9}$$

Under these circumstances, the reward does not incentivize the model to increase or decrease P_1 or P_2 . Essentially, this weighting strategy amplifies the reward for rare classes, enabling the model to learn more effectively from such examples and thereby approximating balanced training.

Models	Methods	Math-500	GSM8k	Minerva	AMC23	Olympiad	Avg
		Generation Accuracy (%)					
Mistral-7B-Instruct-v0.3	w/o RL	10.8	46.7	7.4	5.0	1.8	14.3
	C-RL	19.8	54.2	9.8	12.5	4.2	20.1
	G-RL	47.6	77.8	18.6	30.0	15.6	37.9
	GC-RL	52.1	81.2	19.4	32.5	17.9	40.6
Llama-3.1-8B-Instruct	w/o RL	48.2	84.2	19.1	25.0	15.0	38.3
	C-RL	55.6	87.6	21.3	32.5	20.1	43.4
	G-RL	72.3	92.0	28.9	47.5	28.4	53.8
	GC-RL	75.8	92.6	30.1	50.0	31.6	56.0
		Critique Accuracy (%)					
Mistral-7B-Instruct-v0.3	C-RL	65.2	70.1	60.3	66.4	63.2	65.0
	GC-RL	69.0	74.8	64.5	71.6	66.1	69.2
Llama-3.1-8B-Instruct	C-RL	77.8	80.5	58.3	70.4	67.3	70.9
	GC-RL	81.4	83.2	64.1	75.4	70.6	74.9

Table 4: Generation and critique capabilities evaluation results on Mistral-7B and Llama3.1-8B-Instruct. GC-RL outperforms all other RL training methods in both generation and critique capabilities.

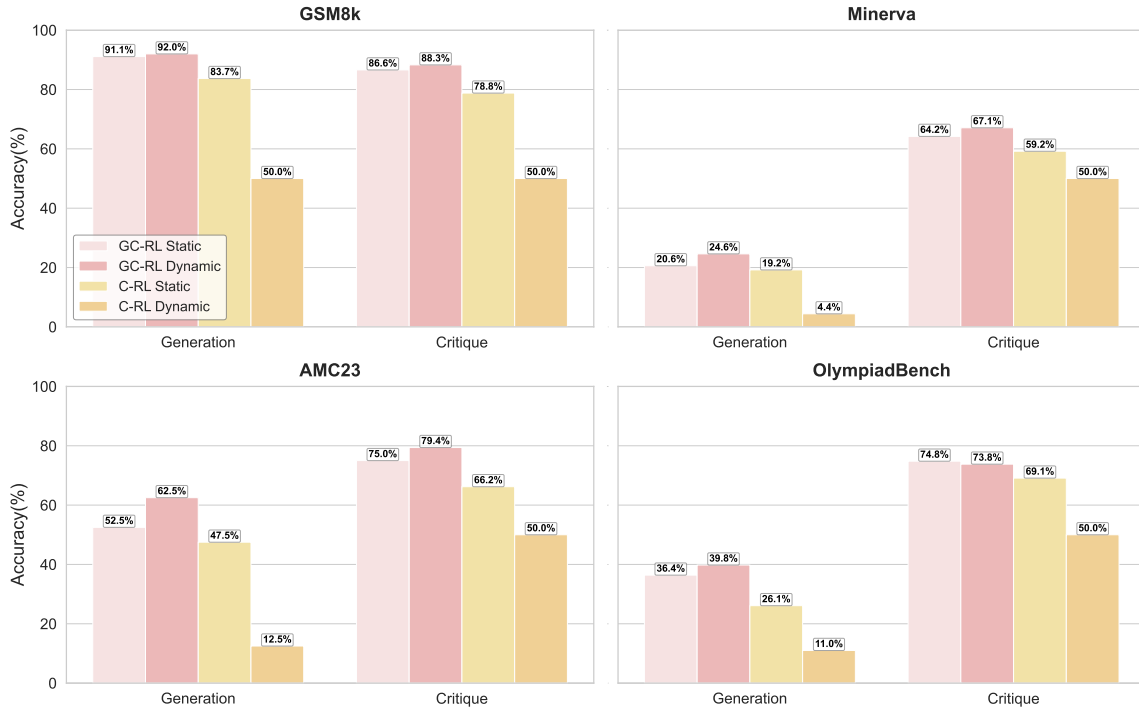


Figure 7: Performance of Qwen2.5-7B on 4 datasets under GC-RL and C-RL settings with static and dynamic critique training data. Dynamic data outperforms static data in the GC-RL setting, while the opposite holds for C-RL.

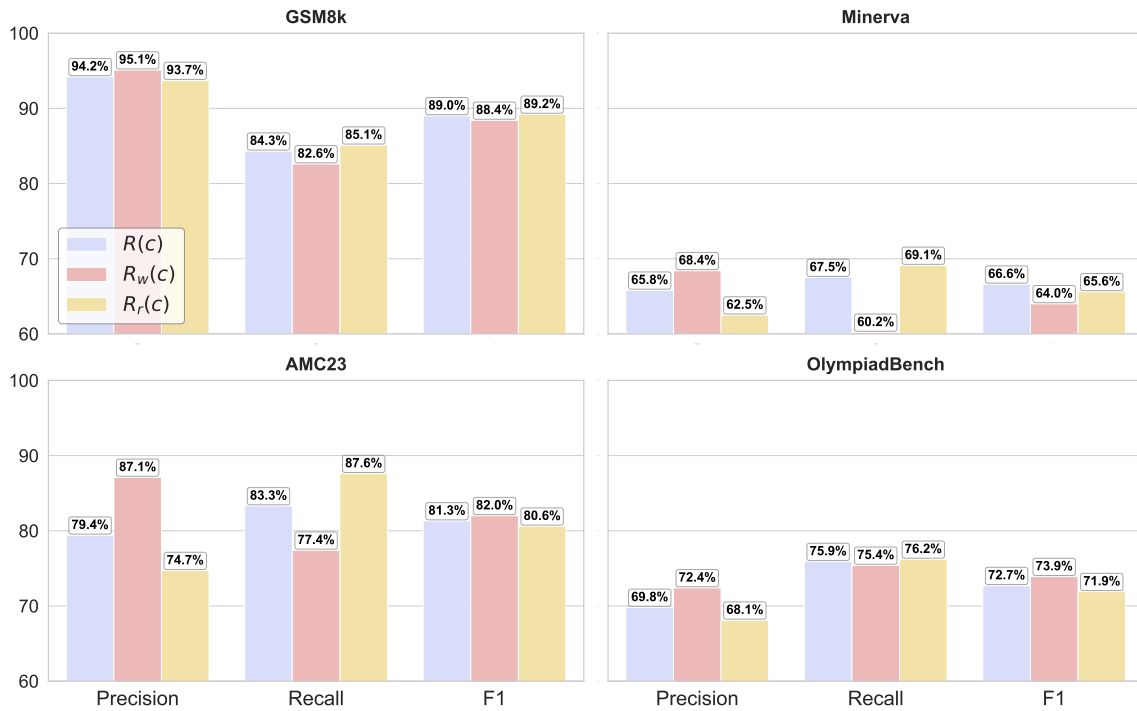


Figure 8: A comparison of critique performance of Qwen2.5-7B with different reward functions on 4 datasets. Compared to baseline $R(c)$, $R_w(c)$ leads to a higher precision while $R_r(c)$ generates a higher recall.

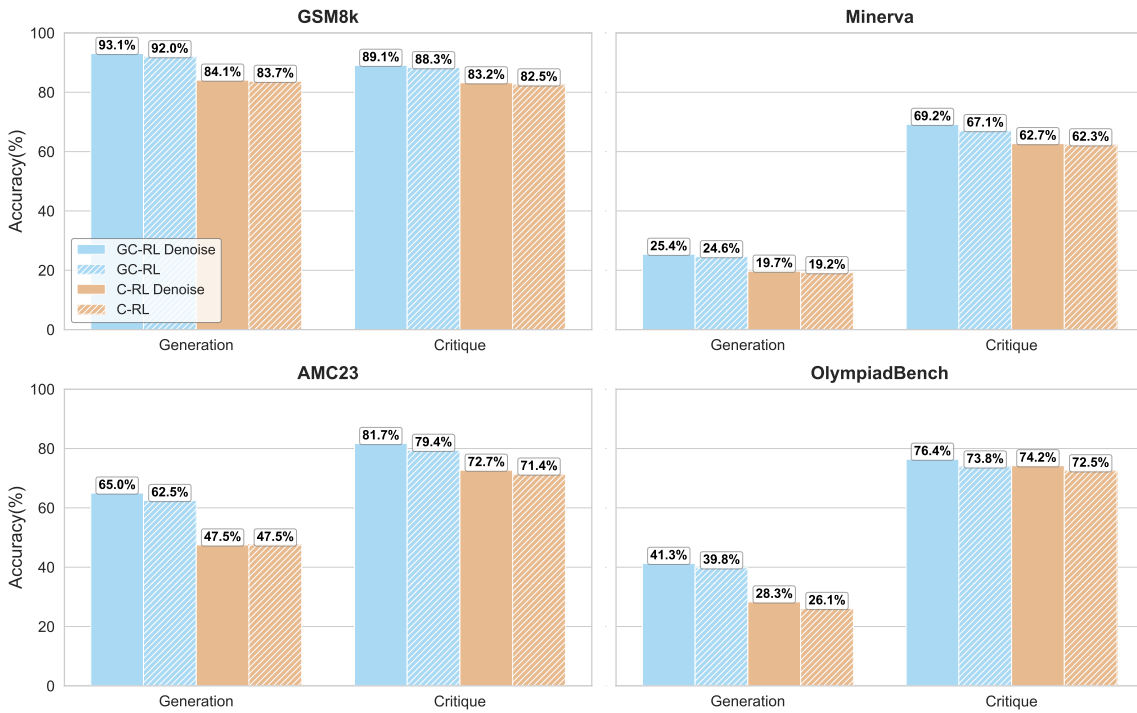


Figure 9: A comparison of model performance of Qwen2.5-7B with/without reward denoising strategy on 4 datasets. In both GC-RL and C-RL settings, reward denoising can improve model performance on both generation and critique capabilities.