

# Dual Consensus: Escaping from Spurious Majority in Unsupervised RLVR via Two-Stage Vote Mechanism

Anonymous ACL submission

## Abstract

Current label-free RLVR approaches for large language models (LLMs), such as TTRL and Self-reward, have demonstrated effectiveness in improving the performance of LLMs on complex reasoning tasks. However, these methods rely heavily on accurate pseudo-label estimation and converge on spurious yet popular answers, thereby trapping in a dominant mode and limiting further improvements. Building on this, we propose **Dual Consensus Reinforcement Learning (DCRL)**, a novel self-supervised training method which is capable of generating more reliable learning signals through a two-stage consensus mechanism. The model initially acts as an *anchor*, producing dominant responses; then it serves as an *explorer*, generating diverse auxiliary signals via a temporary unlearning process. The final training target is derived from the harmonic mean of these two signal sets. Notably, the process operates entirely without external models or supervision. Across eight benchmarks and diverse domains, DCRL consistently improves Pass@1 over majority vote while yielding more stable training dynamics. These results demonstrate that DCRL establishes a scalable path toward stronger reasoning without labels.

## 1 Introduction

Reinforcement Learning with Verifiable Rewards (RLVR) has emerged as an effective approach to boosting the performance of Large Language Models (LLMs), enabling superior reasoning capabilities through long chain-of-thought (Wei et al., 2023) reasoning on various challenging benchmarks (OpenAI et al., 2024; DeepSeek-AI et al., 2025; Yang et al., 2025). However, typically implemented via algorithms such as Group Relative Policy Optimization (GRPO) (Shao et al., 2024), current RLVR approaches heavily rely on human-annotated datasets or, at a minimum, environments that provide verifiable ground-truth signals (Le

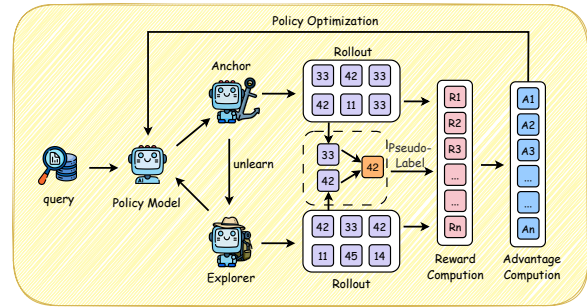


Figure 1: An overview of *Dual Consensus Reinforcement Learning (DCRL)*. Specifically, the policy model assumes two roles: (1) an *anchor* that generates dominant and reliable responses; (2) an *explorer* that produces diverse auxiliary signals through a temporary unlearning process.

et al., 2022; Wang et al., 2024a). This reliance restricts its generalizability to fully unlabeled or distribution-shifted tasks where neither human annotations (Ziegler et al., 2020) nor executable environments are accessible. As LLMs approach or surpass human-level performance, they will inevitably operate in domains where even expert humans cannot provide definitive judgments or reliable evaluations, which motivates the exploration of training on unlabeled data.

Recent studies (Shafayat et al., 2025; Zhao et al., 2025a) find that LLMs can achieve self-improvement without labeled data. One approach is determinism-based methods (Prabhudesai et al., 2025; Zhang et al., 2025b), which derive rewards from the confidence of a single policy along trajectories, thereby encouraging low-entropy and high-confidence predictions. These methods can achieve better performance by sharpening the model’s outputs, but it remains debatable whether they are truly effective for improving reasoning capabilities. Another approach is aggregation-based methods (Zhang et al., 2025c; Yu et al., 2025b; Wu et al., 2025), which derive rewards from agreement across multiple samples, assuming that cross-sample con-

sistency correlates with correctness. Nevertheless, these approaches still suffer from two critical limitations:

- **Spurious Reward Signals:** Models struggle to generate distinguishable reward signals, especially when tackling hard reasoning tasks; the answers derived from majority vote may themselves suffer from systematic biases (Zhao et al., 2025b). In the later stages of training, spurious majority outcomes can come to dominate, yet the correct solutions may instead lie within the minority rollouts.
- **Lack Exploration Capability:** By continuously rewarding consensus across diverse trajectories, the model’s output distribution becomes increasingly rigid and concentrated, resulting in a severe deficiency in exploration capability or even entropy collapse (Cui et al., 2025). Consequently, the model tends to converge to a narrow set of suboptimal responses and exhibits degraded performance when confronted with out-of-domain (OOD) tasks.

In this paper, we propose *Dual Consensus*—a novel framework for Unsupervised Reinforcement Learning with Verifiable Rewards (URLVR) driven by a multi-stage vote mechanism. Our core insight stems from the following intuition: valid reasoning trajectories should not only converge to the dominant mode but also exhibit enhanced robustness when the distribution of the model is artificially flattened.

Instead of naively adopting the fragile majority vote as the pseudo-label, we decompose the rollout process into two stages: *anchor* and *explorer*. The *anchor* stage involves normal rollouts where the model generates responses under its current policy, capturing the dominant reasoning mode. Subsequently, in the *explorer* stage, we introduce a temporary unlearning process to flatten the distribution and enhance exploration, thereby encouraging the generation of diverse auxiliary responses that deviate from the dominant mode. After obtaining the two signal sets, we compute the harmonic mean of their consensus scores to determine the final reward signal. This harmonic mean balances the reliability of the dominant mode (from the anchor stage) and the diversity of potential valid trajectories (from the explorer stage), effectively mitigating the adverse impact of majority vote when it converges to spurious answers.

To validate our approach, we demonstrate the effectiveness of *Dual Consensus* through extensive experiments. We first train models on the large-scale DAPO-14K-Math dataset and evaluate them on multiple established benchmarks. Additionally, we apply test-time adaptation on five distinct datasets to further assess our method. In summary, our key contributions are as follows:

- A novel URLVR method, *Dual Consensus*, which utilizes the intrinsic robustness of the model to guide the model to evolve with policy optimization methods such as GRPO. Notably, the framework is entirely free of external models and enables continuous self-improvement without supervision.
- A pseudo-label selection mechanism and reward design that exploits the intrinsic robustness of the model itself to generate more reliable reward signals, mitigating the toxicity of majority vote when it fails.
- We empirically verify the general effectiveness of *Dual Consensus* in boosting LLMs’ reasoning performance via comprehensive experiments, and additionally present systematic ablation studies and in-depth further analyses.

## 2 Methodology

In this section, we present the details of *Dual Consensus*. It mitigates spurious majority bias by generating diverse signals via Unlearn Then Explore, identifying more accurate labels through Harmonic Election, and stabilizing updates with Adaptive Sampling.

Our framework employs Grouped Relative Policy Optimization (GRPO) (Shao et al., 2024) as its foundational RL algorithm—it stabilizes training by normalizing advantage estimates across multiple rollouts of the same prompt.

For a given input prompt  $x$  (paired with its ground-truth label  $y$ ), GRPO first samples  $n$  rollouts  $\{y_i\}_{i=1}^n$  from the current policy. For each rollout  $y_i$ , GRPO computes a reward  $r_i = R(y_i, y | x)$ , then derives the group-normalized advantage  $\hat{A}_i$ :

$$\hat{A}_i = \frac{r_i - \bar{r}}{\sigma_r} \quad (1)$$

where  $\bar{r} = \frac{1}{n} \sum_{k=1}^n r_k$  denotes the mean reward of the rollout group, and  $\sigma_r$  is the standard deviation of the group rewards. The GRPO objective

optimizes the target policy  $\pi_\theta$  by maximizing the clipped normalized advantage:

$$\mathcal{L}_{\text{GRPO}}(x, y; \theta) = \frac{1}{n} \sum_{i=1}^n \min(\rho_i(\theta) \hat{A}_i, \text{clip}(\rho_i(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_i) - \beta \mathbb{D}_{\text{KL}}[\pi_\theta(\cdot | x) \| \pi_{\text{ref}}(\cdot | x)] \quad (2)$$

where  $\rho_i(\theta) = \pi_\theta(y_i | x) / \pi_{\text{old}}(y_i | x)$  is the importance sampling ratio.

## 2.1 Unlearn Then Explore

Following the GRPO framework, we first initialize an *anchor* model (parameterized by  $\theta'$ ) by cloning the current policy model (parameterized by  $\theta$ ), i.e.,  $\theta' \leftarrow \theta$ . We then strategically apply an unlearning strategy (Liu et al., 2024) to transform this *anchor* into an *explorer*. Unlike EEPO (Chen et al., 2025), which simply employs unlearning as a supplementary technique to enhance exploration, our approach leverages unlearning as a core methodology to actively search for correct answers.

To implement this unlearning strategy, we firstly introduce the standard *negative log-likelihood* (NLL) loss that is defined as:

$$\mathcal{L}_{\text{NLL}} = -\log \pi_{\text{anchor}}(y_{i,t} | x, y_{i,<t}) \quad (3)$$

This loss function imposes heavier penalties on predictions with low probability. Conversely, we take a complementary loss function that inverts the penalty pattern of the NLL loss.

To ensure numerical stability when  $\pi_{\text{anchor}}(y_{i,t} | x, y_{i,<t}) \rightarrow 1$  (which would make  $1 - \pi_{\text{anchor}}(y_{i,t} | x, y_{i,<t})$  approach zero and cause numerical overflow), we first perform a clipping operation on the prediction probability from the *anchor* model:

$$p_{\text{clip}} = \text{clip}(\pi_{\text{anchor}}(y_{i,t} | x, y_{i,<t}), \epsilon, 1 - \epsilon) \quad (4)$$

where  $\epsilon$  is a small positive constant for numerical stability. This clipping operation both avoids the term  $1 - \pi_{\text{anchor}}(y_{i,t} | x, y_{i,<t})$  from becoming excessively small and eliminates unnecessary penalties for tokens with extremely low probabilities that are irrelevant to meaningful exploration.

Based on the clipped probability  $p_{\text{clip}}$ , we define the stabilized unlearning loss as follows:

$$\mathcal{L}_{\text{unlearn}} = -\log(1 - p_{\text{clip}}) \quad (5)$$

This loss function achieves the targeted penalty characteristic: high  $p_{\text{clip}}$  tokens lead to a large negative  $\log(1 - p_{\text{clip}})$ , and minimizing the loss enforces

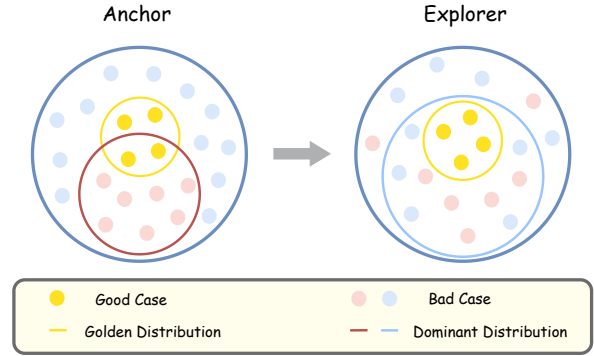


Figure 2: Output distributions of the *Anchor* and *Explorer* models. The *Explorer* model, after the unlearning process, generates a more diverse distribution.

a sharp reduction in their probabilities with strong penalties.

To implement this unlearning process in practice, we apply a single gradient descent step to the *anchor* model using the unlearn loss  $\mathcal{L}_{\text{unlearn}}$ , which transforms it into an *explorer* model:

$$\theta' \leftarrow \theta' - \eta \nabla_{\theta'} \mathcal{L}_{\text{unlearn}}(\theta') \quad (6)$$

where  $\eta$  denotes the learning rate for this unlearning step. Critically, this update is *temporary*—it is confined exclusively to the *anchor* model within the current iteration, and the parameters of the original policy model  $\theta$  remain unchanged.

The gradient update direction  $-\eta \nabla_{\theta'} \mathcal{L}_{\text{unlearn}}$  suppresses high-confidence tokens from the anchor model. The derivative of  $\mathcal{L}_{\text{unlearn}} = -\log(1 - p_{\text{clip}})$  assigns larger gradient magnitudes to tokens with higher  $p_{\text{clip}}$  values, and gradient descent directly drives the reduction of their probabilities to achieve the unlearning effect.

As illustrated in Fig. 2, the explorer model modulates the anchor’s probability distribution, expanding the coverage of generated responses to include trajectories that deviate from the dominant mode, which provides the necessary diverse reasoning signals for subsequent consensus aggregation.

## 2.2 Harmonic Election

Instead of relying on majority vote, we employ the harmonic mean to establish consensus between the anchor and explorer models, which effectively balances exploration and exploitation (shown in Fig 3). This consensus process proceeds in three sequential steps.

First, we perform anchor rollout: drawing  $G$  trajectories from the anchor policy  $\pi_{\text{anchor}}$ , denoted

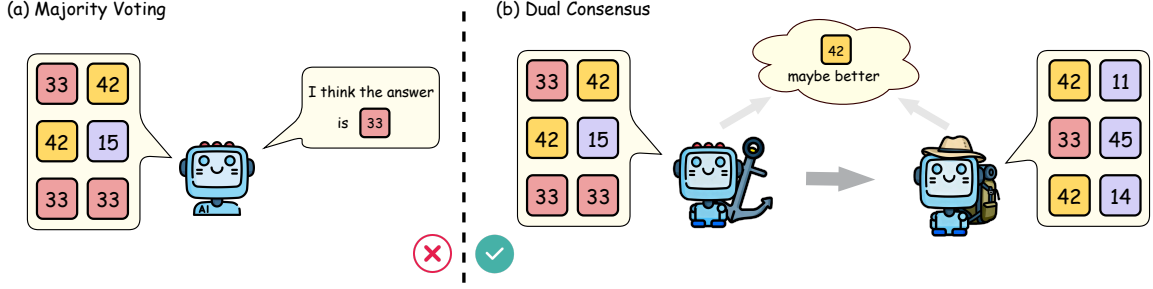


Figure 3: Comparison between *Majority vote* and *Dual Consensus*: Majority vote tends to fall into spurious consensus by over-relying on dominant but potentially incorrect response modes, while *Dual Consensus* mitigates this issue by converting the anchor model (which captures dominant reasoning patterns) into an explorer model via temporary unlearning. This transformation enables the framework to explore diverse alternative response modes, thereby balancing the reliability of current dominant patterns and the diversity of potential valid alternatives, and ultimately achieving more accurate answer selection.

as:

$$O_0 = \{o_1, o_2, \dots, o_G\}. \quad (7)$$

Then we apply temporary unlearning to  $\pi_{\text{anchor}}$  by minimizing the unlearning loss (Eq. 5), thereby suppressing its dominant generation patterns and converting it into the explorer model  $\pi_{\text{explorer}}$ .

We subsequently conduct explorer rollout: generating another  $G$  trajectories from  $\pi_{\text{explorer}}$ , forming the set:

$$O_1 = \{o'_1, o'_2, \dots, o'_G\}. \quad (8)$$

Let  $\mathcal{A}$  be the set of all candidate answers. For each  $a \in \mathcal{A}$ , we compute its empirical occurrence probabilities in  $O_0$  and  $O_1$ :

$$p_0(a) = \frac{1}{G} \sum_{i=1}^G \mathbb{I}(o_i \mapsto a), \quad (9)$$

$$p_1(a) = \frac{1}{G} \sum_{i=1}^G \mathbb{I}(o'_i \mapsto a), \quad (10)$$

where  $\mathbb{I}(\cdot)$  is the indicator function.

The consensus pseudo-label  $y^*$  is then selected as the answer that maximizes the harmonic mean of  $p_0(a)$  and  $p_1(a)$ :

$$y^* = \arg \max_{a \in \mathcal{A}} \frac{2p_0(a)p_1(a)}{p_0(a) + p_1(a)}. \quad (11)$$

We note *Self-Harmony* (Wang et al., 2025) also uses the harmonic mean, yet our method only processes the model’s distribution for the same input, thus avoiding Self-Harmony’s semantic inconsistency in question rephrasing—an frequently recurring issue with catastrophic consequences.

During reward computation, a full reward is assigned to trajectories generating the consensus

pseudo-label  $y^*$ . Trajectories consistent with the majority result of the anchor are assigned a modest reserved reward to avoid negative advantages in subsequent estimation, since such trajectories are still deemed promising relative to other counterparts. Formally, the reward for trajectory  $i$  is defined as:

$$r_i = \begin{cases} 1 & \text{if } y_i(o) = y^*, \\ 0.5 & \text{if } y_i(o) = \hat{y}_{\text{anchor}}, \\ 0 & \text{otherwise,} \end{cases} \quad (12)$$

where  $\hat{y}_{\text{anchor}} = \arg \max_{a \in \mathcal{A}} \sum_{o_i \in O_0} \mathbb{I}(o_i \mapsto a)$  denotes the majority answer from the anchor trajectories.

### 2.3 Adaptive Sampling

Unlike the static sampling strategy in majority vote (Zuo et al., 2025), we introduce the *consensus rate* as a signal to adaptively regulate the contribution of *anchor* and *explorer* rollouts during policy updates.

We formally define the consensus rate  $\rho_t$  at step  $t$  as the proportion of anchor-sampled trajectories whose generated answers are consistent with the majority-voted result  $\hat{y}_{\text{anchor}}$  of the anchor stage, which is formulated as:

$$\rho_t = \frac{1}{|O_0|} \sum_{o \in O_0} \mathbb{I}(y(o) = \hat{y}_{\text{anchor}}) \quad (13)$$

where  $y(o)$  denotes the answer extracted from trajectory  $o$ , and  $\mathbb{I}(\cdot)$  is the indicator function.

To capture long-term consistency and mitigate step-wise noise, we maintain a sliding window of the consensus rate over the most recent  $K$  steps

and compute its mean:

$$\bar{\rho}_t = \frac{1}{K} \sum_{k=1}^K \rho_{t-k+1}. \quad (14)$$

A high  $\bar{\rho}_t$  indicates that the policy model consistently converges to the same answer, reflecting strong model certainty and deterministic behavior; conversely, a low  $\bar{\rho}_t$  suggests ongoing exploration.

Based on  $\bar{\rho}_t$ , we design a dynamic sampling selection rule with threshold  $1/2$ :

- When  $\bar{\rho}_t \leq \frac{1}{2}$ , only anchor trajectories  $O_0$  are used for policy updates. Explorer trajectories  $O_1$  still participate in consensus formation via harmonic vote but are excluded from gradient computation. This prevents premature incorporation of noisy exploratory signals while preserving rare answers that align with the pseudo-label.
- When  $\bar{\rho}_t > \frac{1}{2}$ , both  $O_0$  and  $O_1$  are included in training. This enables the policy to leverage reliable anchor behaviors while actively integrating diverse, high-quality explorations.

The effective rollout set for policy update is defined as:

$$O_{\text{train}} = \begin{cases} O_0 & \text{if } \bar{\rho}_t \leq \frac{1}{2}, \\ O_0 \cup O_1 & \text{if } \bar{\rho}_t > \frac{1}{2}. \end{cases} \quad (15)$$

Only trajectories in  $O_{\text{train}}$  contribute to the policy gradient, ensuring a smooth transition from exploitation-dominant to balanced exploration-exploitation learning.

The overall workflow of the proposed Dual Consensus algorithm is summarized in Algorithm 1.

### 3 Experiments

In this section, we first introduce the experimental setup, then discuss the overall effectiveness of our method, and finally present the results of our ablation studies.

#### 3.1 Setups

To comprehensively evaluate the effectiveness of DCRL, our experiments are conducted under two distinct training paradigms:

- **Large-Scale Unsupervised Learning:** Directly applying DCRL to train models from scratch on a large, unlabeled dataset.
- **Test-Time Adaptation (TTA):** Using DCRL to adapt a pre-trained model to new, unseen benchmarks via constant unsupervised training.

---

#### Algorithm 1 Dual Consensus: An Unsupervised RLVR Algorithm

---

- 1: **Initialize:** policy  $\theta^0$ ; learning rates  $\eta_{\text{GRPO}}, \eta_{\text{u}}$ ; group size  $G$ ; iterations  $T$
  - 2: **for**  $t = 0$  to  $T - 1$  **do**
  - 3:   Sample query  $x \sim \mathcal{D}$
  - 4:   Sample  $G$  trajectories  $O_{\text{anchor}} \sim \pi_{\theta^t}(\cdot | x)$ , compute consensus rate  $\rho_t$  and  $\bar{\rho}_t$
  - 5:    $\theta_e \leftarrow \theta^t - \eta_{\text{u}} \nabla_{\theta^t} \mathcal{L}_{\text{u}}(O_{\text{anchor}})$
  - 6:   Sample  $G$  trajectories  $O_{\text{explorer}} \sim \pi_{\theta_e}(\cdot | x)$
  - 7:    $S(a) = \frac{2p_0(a)p_1(a)}{p_0(a)+p_1(a)}$  for each answer  $a$  in  $O_{\text{anchor}} \cup O_{\text{explorer}}$
  - 8:   Pseudo-label  $y^* = \arg \max_a S(a)$
  - 9:   **if**  $\rho_t > 1/2$  **then**
  - 10:      $O \leftarrow O_{\text{anchor}} \cup O_{\text{explorer}}$
  - 11:   **else**
  - 12:      $O \leftarrow O_{\text{anchor}}$
  - 13:   **end if**
  - 14:   Compute rewards and advantages for  $O$
  - 15:    $\theta^{t+1} \leftarrow \theta^t + \eta_{\text{GRPO}} \nabla_{\theta} J_{\text{GRPO}}(\theta^t; O)$
  - 16: **end for**
- 

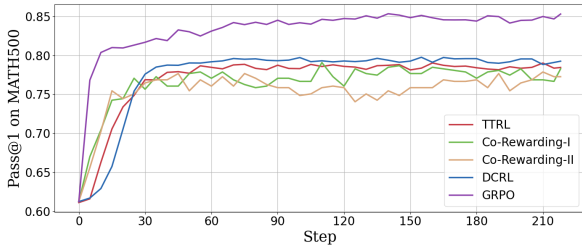
**Models:** Our target models include Llama3.2-3B-Instruct (Grattafiori et al., 2024), Qwen3-4B-Base, and Qwen3-8B-Base (Yang et al., 2025). Additionally, for the TTA paradigm, we also conduct experiments on Qwen2.5-Math-1.5B (Yang et al., 2024).

**Implementation:** Our primary training dataset is DAPO-Math-14k, a processed version of DAPO-Math-17k (Yu et al., 2025a), which is refined by deduplicating prompts and standardizing the formatting of both prompts and reference answers. We train each model on this dataset for two epochs to mitigate overfitting. All experiments are based on the VeRL framework (Sheng et al., 2025). Implementation details are reported in Appendix A.

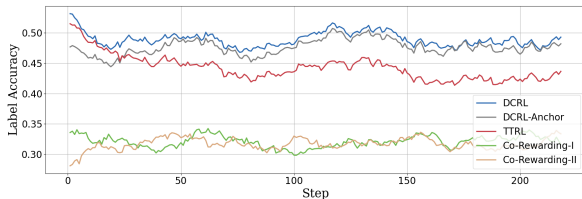
**Benchmarks:** Our benchmark suite comprises eight challenging datasets, including six math-specific: (1) MATH-500 (Hendrycks et al., 2021), (2) GSM8K (Cobbe et al., 2021), (3) AIME24 (Yang et al., 2025), (4) Minerva-math (Lewkowycz et al., 2022), (5) AMC (Zuo et al., 2025), (6) OlympiadBench (He et al., 2024); and two multi-task benchmarks: (1) MMLU-Pro (Wang et al., 2024b), (2) GPQA-Diamond (Rein et al., 2024). For the TTA experiments, we evaluate on subsets of MATH-500, AIME24, AMC, and GPQA-Diamond.

Methods	Math				Multi-Task				Average
	Math	GSM8K	AIME24	Minerva.	AMC	Olympiad.	MMLU.	GPQA.	
<i>Llama3.2-3B-Instruct</i>									
Vanilla	42.4	76.1	4.5	11.7	20.6	14.6	26.4	22.3	27.3
GRPO	49.2	79.7	13.5	14.9	23.2	15.3	32.8	23.8	31.5
RENT	45.4	78.5	9.4	11.2	20.9	15.2	30.3	22.7	29.2
TTRL	44.6	71.8	10.2	12.4	21.5	14.6	25.9	21.4	27.8
Co-Rewarding-I	45.3	77.4	<b>11.8</b>	15.1	22.2	14.8	<b>33.9</b>	24.3	30.6
Co-Rewarding-II	47.6	76.9	11.6	13.7	<b>23.3</b>	15.3	24.9	21.4	29.3
<b>DCRL(ours)</b>	<b>47.8</b>	<b>79.1</b>	<b>11.8</b>	<b>15.8</b>	<b>23.3</b>	<b>15.4</b>	<b>33.9</b>	<b>24.7</b>	<b>31.4</b>
<i>Qwen3-4B-Base</i>									
Vanilla	47.4	87.0	9.3	24.5	38.3	34.9	48.6	30.7	40.0
GRPO	76.8	91.8	13.1	32.9	45.2	36.1	52.2	34.8	47.8
RENT	72.1	85.0	10.2	23.7	44.3	36.6	48.6	32.2	44.0
TTRL	74.4	91.3	11.5	29.0	44.7	37.0	51.2	32.8	46.4
Co-Rewarding-I	<b>74.8</b>	91.6	11.8	31.1	43.1	36.9	52.2	32.9	46.7
Co-Rewarding-II	74.5	91.7	11.6	30.8	44.5	37.1	<b>52.6</b>	33.3	47.0
<b>DCRL(ours)</b>	<b>74.8</b>	<b>91.7</b>	<b>12.4</b>	<b>31.2</b>	<b>45.6</b>	<b>37.2</b>	<b>52.6</b>	<b>34.3</b>	<b>47.4</b>
<i>Qwen3-8B-Base</i>									
Vanilla	61.8	88.3	12.5	25.0	48.9	40.1	50.6	32.3	44.9
GRPO	82.0	93.8	20.4	34.1	57.0	45.8	57.5	38.1	53.5
RENT	75.3	86.5	13.1	25.4	49.7	40.3	52.5	34.5	47.1
TTRL	78.3	92.9	14.4	<b>32.7</b>	51.2	<b>40.9</b>	55.2	36.5	50.2
Co-Rewarding-I	78.2	92.7	12.2	32.0	51.2	40.7	<b>56.7</b>	<b>37.9</b>	50.2
Co-Rewarding-II	77.2	92.9	12.2	32.3	51.8	40.2	56.1	37.6	50.0
<b>DCRL(ours)</b>	<b>79.2</b>	<b>93.3</b>	<b>14.7</b>	<b>32.7</b>	<b>51.9</b>	<b>40.9</b>	<b>56.7</b>	<b>37.9</b>	<b>50.9</b>

Table 1: **Main Results (%) of DCRL Trained on DAPO-Math-14k: DCRL Outperforms All Label-Free Baselines.** The best results are highlighted in **bold**. DCRL exceeds all unsupervised methods and nearly matches GRPO with gold labels.



(a) Accuracy curve on the MATH500 benchmark.



(b) Label accuracy curve (smoothed).

Figure 4: Training Dynamics of Dual Consensus on Qwen3-8B-Base. DCRL-Anchor in Fig. 4b refers to the majority vote of the anchor model in DCRL.

**Baselines:** We compare DCRL against four unsupervised RLVR methods, including one determinism-based approach **RENT** (Prabhudesai et al., 2025), and three aggregation-based ap-

proaches **TTRL** (Zuo et al., 2025), **Co-Rewarding-I**, and **Co-Rewarding-II** (Zhang et al., 2025c). Specifically, for Co-Rewarding-I, we adopt a dataset rephrased by Qwen3-32B for all related experiments.

**Metrics:** We use the Pass@1 metric. For each question, we sample 16 predictions using a temperature of 0.6 and a top-p value of 0.95. The final reported score is the average Pass@1 accuracy across these 16 independent seeds.

## 3.2 Results

### 3.2.1 Main Performance of Dual Consensus

Table 1 presents the main results of DCRL trained on DAPO-Math-14k across eight challenging reasoning benchmarks. Our method consistently outperforms all label-free baselines—including RENT, TTRL, and both variants of Co-Rewarding—across different model scales and task domains. Notably, on the Qwen3-8B-Base model, DCRL achieves 79.2% on MATH-500, surpassing the strongest baseline TTRL at 78.3% by 0.9%, and improves AIME24 from 14.4% to 14.7%, demonstrating

its effectiveness on extremely hard competition-level problems. On multi-task benchmarks, DCRL matches or exceeds Co-Rewarding-I on MMLU-Pro and GPQA-Diamond, confirming its generalizability beyond pure math reasoning.

Remarkably, despite being fully unsupervised and using no ground-truth labels, DCRL achieves performance on par with, and occasionally exceeds, the supervised GRPO baseline. On Llama3.2-3B-Instruct, Qwen3-4B-Base, and Qwen3-8B-Base, it yields average gains of +7.5%, +4.0%, and +6.1% on MMLU-Pro respectively.

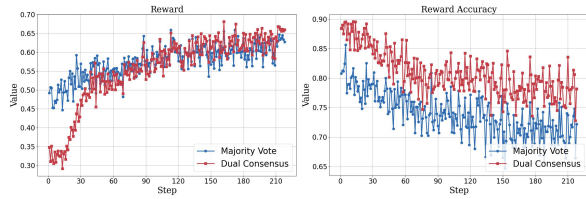


Figure 5: Comparison of reward signal between Majority Vote and Dual Consensus.

**Compared with Majority Vote:** We adopt TTRL (Zuo et al., 2025) as a representative baseline that relies on standard majority vote for pseudo-label estimation. Fig 4 and Fig 5 depict the training dynamics on Qwen3-8B-Base: our method can select low-consistency answers (especially early in training) and achieves higher ground-truth reward accuracy. This demonstrates that Dual Consensus produces more reliable supervision signals than naive majority vote.

### 3.2.2 Test-time Adaption

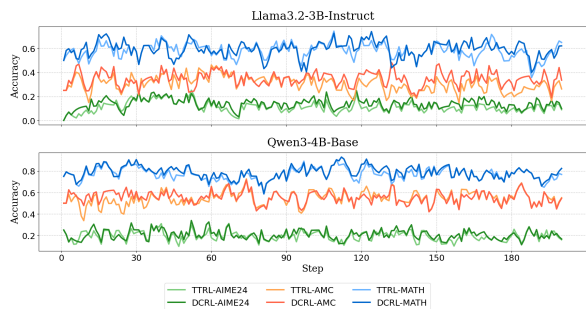


Figure 6: Label accuracy curves (smoothed) of test-time adaption for DCRL and TTRL across different tasks.

Test-time Adaption (TTA) serves as a critical validation scenario for unsupervised RLVR methods, as it directly evaluates the ability to escape spurious majority bias and generalize to unseen reasoning tasks without labeled data. The performance

Methods	Datasets				Average
	MATH500	AIME24	AMC	GPQA	
<i>Qwen2.5-Math-1.5B</i>					
Vanilla	30.6	5.6	23.4	15.4	18.7
TTRL	72.9	17.0	45.3	22.2	39.3
<b>DCRL(ours)</b>	74.5	17.7	46.6	22.8	40.4
$\Delta(-TTRL)$	$\uparrow 2.2\%$	$\uparrow 4.1\%$	$\uparrow 2.8\%$	$\uparrow 2.7\%$	$\uparrow 2.7\%$
<i>Llama3.2-3B-Instruct</i>					
Vanilla	42.4	4.5	20.6	22.3	22.4
TTRL	59.6	10.0	26.5	28.9	31.2
<b>DCRL(ours)</b>	59.8	13.3	32.3	32.6	34.5
$\Delta(-TTRL)$	$\uparrow 0.3\%$	$\uparrow 33.0\%$	$\uparrow 21.8\%$	$\uparrow 12.8\%$	$\uparrow 10.5\%$
<i>Qwen3-4B-Base</i>					
Vanilla	47.4	9.3	38.3	30.7	31.4
TTRL	82.6	17.2	56.5	35.3	47.9
<b>DCRL(ours)</b>	83.4	20.6	56.5	35.6	49.0
$\Delta(-TTRL)$	$\uparrow 0.9\%$	$\uparrow 19.7\%$	$\uparrow 0.0\%$	$\uparrow 0.8\%$	$\uparrow 2.2\%$
<i>Qwen3-8B-Base</i>					
Vanilla	61.8	12.5	48.9	32.3	38.8
TTRL	85.7	19.8	59.8	43.1	52.1
<b>DCRL(ours)</b>	86.4	22.8	61.0	44.5	53.6
$\Delta(-TTRL)$	$\uparrow 0.8\%$	$\uparrow 15.1\%$	$\uparrow 2.0\%$	$\uparrow 3.2\%$	$\uparrow 2.8\%$

Table 2: **Results (%) of Test-Time Adaption Trained on Different Datasets: DCRL Consistently Outperforms the TTRL Baseline.** All results are evaluated under the same experimental settings and reported as the average pass@1 over 16 independent seeds.

of TTA is presented in Table 2. We did not include methods such as *RESTRAIN* (Yu et al., 2025b) and *Self-Harmony* (Liu et al., 2025) as baselines because their official code has not been released.

Although TTA still enables generalization to other unseen scenarios with limited data (Shafayat et al., 2025; Zuo et al., 2025), a defining characteristic of TTA is its immunity to overfitting concerns. In this setting, the model’s ability to a priority avoid spurious reward signals becomes critically important. As demonstrated in Table 2, our DCRL consistently outperforms TTRL across all evaluated tasks, validating the efficacy of our dual consensus mechanism in suppressing misleading signals.

### 3.3 Ablation Studies

To understand the contribution of each component in our DCRL framework, we conduct a series of ablation studies on the Qwen3-8B-Base model, and the results are summarized in Table 3. More detailed results are shown in Appendix A.

**Impact of Harmonic Election** Replacing harmonic mean consensus with simple majority voting from the anchor model alone leads to performance drops. This confirms that harmonic election effectively mitigates spurious majority bias by fusing

Method	MATH	GSM8K	AIME24	MMLU.	GPQA.	Average
<b>DCRL (Full)</b>	<b>79.2</b>	93.3	<b>14.7</b>	56.7	<b>37.9</b>	<b>56.3</b>
- w/o Harmonic Election	78.7	<b>93.4</b>	14.3	56.2	37.8	56.0
- w/o Conservative Reward	78.3	93.3	12.2	<b>57.2</b>	36.2	55.4
- w/o Dynamic Sampling	76.4	90.3	14.3	52.0	34.4	53.4

Table 3: Ablation Studies to Analyze the Contribution of DCRL Core Modules with the Qwen3-8B-Base Model.

both dominant and diverse exploratory signals to produce more reliable pseudo-labels.

**Impact of Conservative Reward** Simplifying our reward design to a binary scheme (1 for correct, 0 otherwise) results in performance degradation, especially in difficult tasks. This demonstrates that our conservative reward, which reserves a modest reward for anchor majority answers, stabilizes training by preventing extreme fluctuations in advantage estimation and avoiding harsh penalties to high-confidence trajectories.

**Impact of Dynamic Sampling** Using both anchor and explorer samples for training at all times leads to the worst overall performance. This underscores the importance of dynamic sampling in balancing exploration and exploitation: it excludes noisy signals early on to avoid reward hacking and incorporates high-quality exploration later, ensuring stable training while preserving the ability to escape suboptimal modes.

## 4 Related Works

**Unsupervised RL for LLMs:** LLMs can achieve self-improvement without labeled data via two typical unsupervised RL paradigms: determinism-based methods (Prabhudesai et al., 2025; Zhang et al., 2025b) encourage low-entropy and high-confidence predictions for performance sharpening, while aggregation-based methods (Zhang et al., 2025c; Yu et al., 2025b; Wu et al., 2025) assign rewards by cross-sample agreement, which takes cross-sample consistency as the proxy of prediction correctness.

**Test-time Adaptation for LLMs:** Recent works (Akyürek et al., 2025) have demonstrated that LLMs can leverage reinforcement learning to conduct test-time adaptation (Sun et al., 2020), which effectively enhances model performance on unseen data and even surpasses the performance of standard training protocols. This paradigm (Zuo et al., 2025; Wu et al., 2025; Liu et al., 2025; Zhou et al.,

2025; Zhang et al., 2025a) empowers models to dynamically adapt to novel task distributions without access to extra labeled training data.

## 5 Conclusion

In this paper, we propose DCRL, an unsupervised Reinforcement Learning with Verifiable Rewards (RLVR) framework that transforms intrinsic model robustness into reliable learning signals, enabling LLMs to self-improve on reasoning tasks without annotated data. By (i) adopting an Unlearn-Then-Explore strategy to break dominant suboptimal reasoning patterns and enhance exploration capability, (ii) leveraging a Harmonic Election mechanism to balance reliability and diversity for robust pseudo-label estimation, and (iii) introducing Adaptive Sampling to dynamically regulate the exploration-exploitation trade-off during training, DCRL effectively mitigates the spurious majority bias—a critical limitation of existing label-free RLVR methods. Empirically, extensive evaluations across diverse LLMs and challenging reasoning benchmarks demonstrate that DCRL consistently outperforms current determinism-based approaches and aggregation-based approaches, which paves a scalable path for LLM self-improvement without external supervision.

## Limitations

Although DCRL successfully mitigates spurious majority issues and boosts reasoning performance via dual consensus and enhanced exploration, it still has key limitations when encountering severe systematic prior bias, where both anchor and explorer signals converge to consistent spurious consensus. This provides little corrective supervision and even reinforces misleading reasoning patterns through policy optimization. Moreover, its performance gains diminish for extremely complex out-of-distribution reasoning tasks that deviate far from the model’s pretraining distribution, as *anchor-explorer* fails to reconstruct novel reasoning paths and dual consensus signals become unreliable.

## References

Ekin Akyürek, Mehul Damani, Adam Zweiger, Linlu Qiu, Han Guo, Jyothish Pari, Yoon Kim, and Jacob Andreas. 2025. [The surprising effectiveness of test-time training for few-shot learning](#). *Preprint*, arXiv:2411.07279.

Liang Chen, Xueting Han, Qizhou Wang, Bo Han, Jing Bai, Hinrich Schutze, and Kam-Fai Wong. 2025. [Eepo: Exploration-enhanced policy optimization via sample-then-forget](#). *Preprint*, arXiv:2510.05837.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *Preprint*, arXiv:2110.14168.

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

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *Preprint*, arXiv:2501.12948.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.

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 problems](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3828–3850, Bangkok, Thailand. Association for Computational Linguistics.

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](#). *Preprint*, arXiv:2103.03874.

Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. 2022. [Large language models can self-improve](#). *Preprint*, arXiv:2210.11610.

Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven C. H. Hoi. 2022. [Coder1: Mastering code generation through pretrained models and deep reinforcement learning](#). *Preprint*, arXiv:2207.01780. 587  
588  
589  
590  
591

Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. [Solving quantitative reasoning problems with language models](#). *Preprint*, arXiv:2206.14858. 592  
593  
594  
595  
596  
597  
598

Jia Liu, ChangYi He, YingQiao Lin, MingMin Yang, FeiYang Shen, and ShaoGuo Liu. 2025. [Ettl: Balancing exploration and exploitation in llm test-time reinforcement learning via entropy mechanism](#). *Preprint*, arXiv:2508.11356. 599  
600  
601  
602  
603

Sijia Liu, Yuanshun Yao, Jinghan Jia, Stephen Casper, Nathalie Baracaldo, Peter Hase, Yuguang Yao, Chris Yuhao Liu, Xiaojun Xu, Hang Li, Kush R. Varshney, Mohit Bansal, Sanmi Koyejo, and Yang Liu. 2024. [Rethinking machine unlearning for large language models](#). *Preprint*, arXiv:2402.08787. 604  
605  
606  
607  
608  
609

OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, and 244 others. 2024. [Openai o1 system card](#). *Preprint*, arXiv:2412.16720. 610  
611  
612  
613  
614  
615  
616

Mihir Prabhudesai, Lili Chen, Alex Ippoliti, Katerina Fragkiadaki, Hao Liu, and Deepak Pathak. 2025. [Maximizing confidence alone improves reasoning](#). *Preprint*, arXiv:2505.22660. 617  
618  
619  
620

David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2024. [GPQA: A graduate-level google-proof q&a benchmark](#). In *First Conference on Language Modeling*. 621  
622  
623  
624  
625

Sheikh Shafayat, Fahim Tajwar, Ruslan Salakhutdinov, Jeff Schneider, and Andrea Zanette. 2025. [Can large reasoning models self-train?](#) *Preprint*, arXiv:2505.21444. 626  
627  
628  
629

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. [Deepseekmath: Pushing the limits of mathematical reasoning in open language models](#). *Preprint*, arXiv:2402.03300. 630  
631  
632  
633  
634  
635

Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2025. [Hybridflow: A flexible and efficient rlhf framework](#). In *Proceedings of the Twentieth European Conference on Computer Systems*, pages 1279–1297. 636  
637  
638  
639  
640  
641

642	Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei A. Efros, and Moritz Hardt. 2020. <a href="#">Test-time training for out-of-distribution generalization</a> .	Xu. 2025b. <a href="#">Restrained: From spurious votes to signals – self-driven rl with self-penalization</a> . <i>Preprint</i> , arXiv:2510.02172.	699
643			700
644			701
645	Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang Sui. 2024a. <a href="#">Math-shepherd: Verify and reinforce llms step-by-step without human annotations</a> . <i>Preprint</i> , arXiv:2312.08935.	Haoyu Zhang, Jiaxian Guo, Yusuke Iwasawa, and Yutaka Matsuo. 2025a. <a href="#">Aqa-ttrl: Self-adaptation in audio question answering with test-time reinforcement learning</a> . <i>Preprint</i> , arXiv:2510.05478.	702
646			703
647			704
648			705
649			
650	Ru Wang, Wei Huang, Qi Cao, Yusuke Iwasawa, Yutaka Matsuo, and Jiaxian Guo. 2025. <a href="#">Self-harmony: Learning to harmonize self-supervision and self-play in test-time reinforcement learning</a> . <i>Preprint</i> , arXiv:2511.01191.	Qingyang Zhang, Haitao Wu, Changqing Zhang, Peilin Zhao, and Yatao Bian. 2025b. <a href="#">Right question is already half the answer: Fully unsupervised llm reasoning incentivization</a> . <i>Preprint</i> , arXiv:2504.05812.	706
651			707
652			708
653			709
654			
655	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhui Chen. 2024b. <a href="#">Mmlu-pro: A more robust and challenging multi-task language understanding benchmark</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 37, pages 95266–95290. Curran Associates, Inc.	Zizhuo Zhang, Jianing Zhu, Xinmu Ge, Zihua Zhao, Zhanke Zhou, Xuan Li, Xiao Feng, Jiangchao Yao, and Bo Han. 2025c. <a href="#">Co-rewarding: Stable self-supervised rl for eliciting reasoning in large language models</a> . <i>Preprint</i> , arXiv:2508.00410.	710
656			711
657			712
658			713
659			714
660			
661		Andrew Zhao, Yiran Wu, Yang Yue, Tong Wu, Quentin Xu, Yang Yue, Matthieu Lin, Shenzi Wang, Qingyun Wu, Zilong Zheng, and Gao Huang. 2025a. <a href="#">Absolute zero: Reinforced self-play reasoning with zero data</a> . <i>Preprint</i> , arXiv:2505.03335.	715
662			716
663			717
664			718
665	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. <a href="#">Chain-of-thought prompting elicits reasoning in large language models</a> . <i>Preprint</i> , arXiv:2201.11903.	Wenting Zhao, Pranjal Aggarwal, Swarnadeep Saha, Asli Celikyilmaz, Jason Weston, and Iliia Kulikov. 2025b. <a href="#">The majority is not always right: RL training for solution aggregation</a> . <i>Preprint</i> , arXiv:2509.06870.	720
666			721
667			722
668			723
669			724
670	Jianghao Wu, Yasmeen George, Jin Ye, Yicheng Wu, Daniel F. Schmidt, and Jianfei Cai. 2025. <a href="#">Spine: Token-selective test-time reinforcement learning with entropy-band regularization</a> . <i>Preprint</i> , arXiv:2511.17938.	Yujun Zhou, Zhenwen Liang, Haolin Liu, Wenhao Yu, Kishan Panaganti, Linfeng Song, Dian Yu, Xiangliang Zhang, Haitao Mi, and Dong Yu. 2025. <a href="#">Evolving language models without labels: Majority drives selection, novelty promotes variation</a> . <i>Preprint</i> , arXiv:2509.15194.	725
671			726
672			727
673			728
674			729
675			730
676	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. <a href="#">Qwen3 technical report</a> . <i>Preprint</i> , arXiv:2505.09388.	Xinyu Zhu, Mengzhou Xia, Zhepei Wei, Wei-Lin Chen, Danqi Chen, and Yu Meng. 2025. <a href="#">The surprising effectiveness of negative reinforcement in llm reasoning</a> . <i>Preprint</i> , arXiv:2506.01347.	731
677			732
678			733
679			734
680			
681	An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. 2024. <a href="#">Qwen2.5-math technical report: Toward mathematical expert model via self-improvement</a> . <i>Preprint</i> , arXiv:2409.12122.	Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2020. <a href="#">Fine-tuning language models from human preferences</a> . <i>Preprint</i> , arXiv:1909.08593.	735
682			736
683			737
684			738
685			739
686			
687			
688	Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, and 16 others. 2025a. <a href="#">Dapo: An open-source llm reinforcement learning system at scale</a> . <i>Preprint</i> , arXiv:2503.14476.	Yuxin Zuo, Kaiyan Zhang, Li Sheng, Shang Qu, Ganqu Cui, Xuekai Zhu, Haozhan Li, Yuchen Zhang, Xinwei Long, Ermo Hua, Biqing Qi, Youbang Sun, Zhiyuan Ma, Lifan Yuan, Ning Ding, and Bowen Zhou. 2025. <a href="#">Ttrl: Test-time reinforcement learning</a> . <i>Preprint</i> , arXiv:2504.16084.	740
689			741
690			742
691			743
692			744
693			745
694			
695			
696	Zhaoning Yu, Will Su, Leitian Tao, Haozhu Wang, Aashu Singh, Hanchao Yu, Jianyu Wang, Hongyang Gao, Weizhe Yuan, Jason Weston, Ping Yu, and Jing	<b>A Implementation Details</b>	746
697		<b>A.1 Prompt</b>	747
698		We use the same suffix prompt both in the training and evaluation of our experiments to promote clear and step-by-step reasoning:	748
			749
			750
			751

Please reason step by step, and put your final answer within `\boxed{}`.

## A.2 Hyperparameters

Hyperparameter settings of our experiment on Qwen3-8B-Base is shown in Table 4.

Table 4: Hyperparameter Settings for DCRL Framework.

Hyperparameter	Value
Batch size	128
Mini batch size	128
Micro batch size	8
Max prompt length	4096
Max response length	3072
Learning rate	$1 \times 10^{-6}$
LR warmup steps ratio	0.1
Learning rate warmup	cosine
Optimizer	Adam
Temperature	0.6
Top k	-1
Top p	0.95
Unlearn LR	$3 \times 10^{-7}$
Number of anchor samples per example	16
Number of explorer samples per example	16
Number of samples per example for policy update	16
Number of samples per example for Testing	16
Use KL loss	False

## A.3 Baseline Implementation

For all baselines, we use the official code provided in their public repositories. For TTRL, we set the learning rate to  $1 \times 10^{-6}$  and the warm-up ratio to 0.1 for large-scale unsupervised learning. For Co-Rewarding-I, we adopt the DAPO-Math-14k dataset rephrased by Qwen3-32B as provided in the original source code. Besides, no external models are used in all baseline experiments. All other hyperparameter settings for the baseline are kept identical to the original configuration.

## B Detailed Results

### B.1 Detailed Results of Ablation Studies

Detailed results of ablation studies are shown in fig 7.

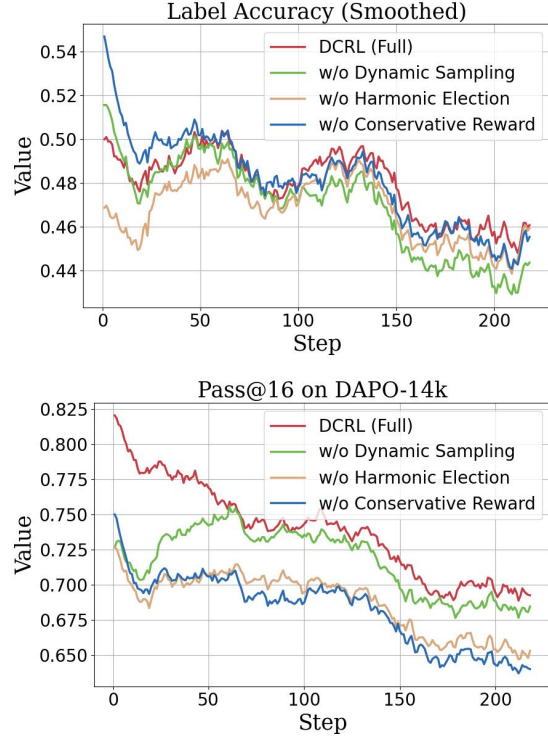


Figure 7: Detailed Results of Ablation Studies with Qwen3-8B-Base, including pass@16 on training dataset and label accuracy

## C Extra Experiments

### C.1 Hyperparameter Sensitivity

A key empirical finding from our experiments is that differing models necessitate tailored unlearning learning rates ( $\eta_u$ ). A value that is too large risks disrupting the model’s ability to generate valid reasoning trajectories, while one that is too small fails to effectively suppress spurious dominant modes. In this section, we present the detailed performance of the Qwen3-8B-Base model under various unlearning learning rate configurations, thereby validating the robustness of our Unlearn-Then-Explore strategy and the rationale behind our specific hyperparameter selection. Results are shown in Table 5.

### C.2 Comparison of Different Consensus Strategies

To further validate the effectiveness of our proposed Harmonic Election mechanism, we compare

Unlearn LR ( $\eta_u$ )	MATH	GSM8K	AIME24	MMLU.	GPQA.	Average
$1 \times 10^{-7}$	78.6	92.7	12.6	55.5	37.2	55.3
$3 \times 10^{-7}$	<b>79.2</b>	<b>93.3</b>	<b>14.7</b>	<b>56.7</b>	<b>37.9</b>	<b>56.3</b>
$5 \times 10^{-7}$	77.6	92.7	11.5	54.3	36.4	54.5
$1 \times 10^{-6}$ (Failed)	67.3	88.1	12.1	50.2	35.1	50.5

Table 5: Sensitivity Analysis of Unlearning Learning Rate ( $\eta_u$ ) on Qwen3-8B-Base.

it against several alternative consensus strategies. These strategies differ in how they aggregate the signals from the *anchor* and *explorer* models to select the final pseudo-label  $y^*$ . The key insight is that a valid reasoning path should be robust, i.e., supported by both the dominant mode (anchor) and the exploratory distribution (explorer).

We evaluate the following strategies:

- **Majority Vote (Anchor Only):** Simply select the majority answer from the *anchor* model’s rollouts.
- **Majority Vote (Anchor + Explorer):** A simple aggregation strategy that combines all rollouts from both the *anchor* and *explorer* models and selects the majority answer.
- **Harmonic Mean (Ours):** Our proposed strategy, which selects the answer that maximizes the harmonic mean of its probabilities in the *anchor* and *explorer* distributions.

The results are presented in Table 6 and Fig 8. We observe that simply combining all samples (*Anchor + Explorer*) does not improve performance and can even be detrimental, as it does not effectively filter out spurious signals. In contrast, our harmonic mean strategy achieves the best overall performance, demonstrating its superior ability to balance reliability and diversity in pseudo-label selection.

## D Why Does Dual Consensus Work?

We formally prove that the dual consensus pseudo-label selection mechanism achieves higher accuracy than naive majority vote by mitigating spurious majority bias, under mild and realistic assumptions.

### D.1 Problem Setup & Definitions

Let  $\mathcal{A}$  be the set of candidate answers, and  $y_{\text{true}} \in \mathcal{A}$  be the ground-truth answer.

- $\pi_{\text{anchor}}(a)$ : Probability of answer  $a$  from the anchor model.

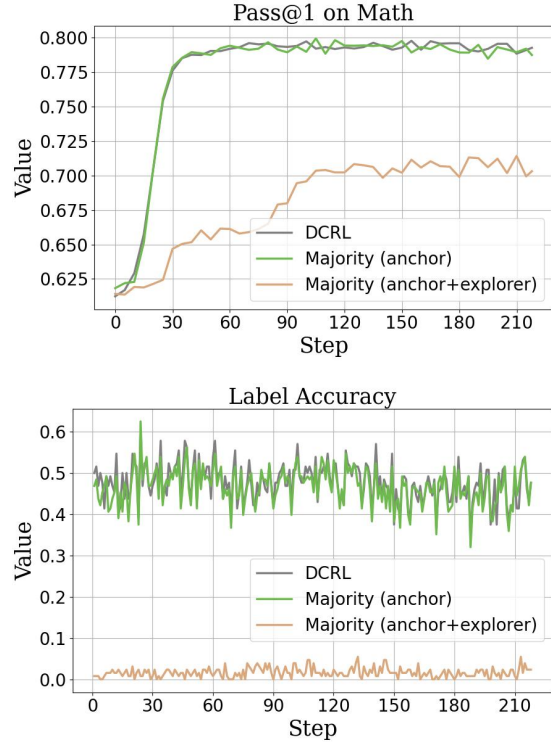


Figure 8: Curves of Different Consensus Strategies for Pseudo-Label Selection on Qwen3-8B-Base.

- $\pi_{\text{explorer}}(a)$ : Probability of answer  $a$  from the explorer model (after unlearning).
- $p_0(a), p_1(a)$ : Empirical probabilities of  $a$  from  $G$  rollouts of the anchor and explorer models, respectively.
- **Majority Vote:**  $\hat{y}_{\text{MV}} = \arg \max_a p_0(a)$ .
- **Dual Consensus:**  $y_{\text{DC}}^* = \arg \max_a S(a)$ , where  $S(a) = \frac{2p_0(a)p_1(a)}{p_0(a)+p_1(a)}$  is the harmonic mean score.

### D.2 Key Assumptions

We introduce three realistic assumptions for LLMs with spurious majority bias.

**Assumption 1 (Spurious Majority Bias):** There exists a spurious dominant answer  $y_{\text{sp}} \neq y_{\text{true}}$  such that:

$$\pi_{\text{anchor}}(y_{\text{sp}}) \gg \pi_{\text{anchor}}(y_{\text{true}})$$

This is the core failure mode of majority vote.

Consensus Strategy	MATH	GSM8K	AIME24	MMLU.	GPQA.	Avg.
Majority Vote (Anchor Only)	78.7	<b>93.4</b>	14.3	56.2	37.8	56.0
Majority Vote (Anchor + Explorer)	64.7	90.3	14.3	52.0	34.4	51.5
Harmonic Mean (Ours)	<b>79.2</b>	93.3	<b>14.7</b>	<b>56.7</b>	<b>37.9</b>	<b>56.3</b>

Table 6: Comparison of Different Consensus Strategies for Pseudo-Label Selection on Qwen3-8B-Base.

**Assumption 2 (Effective Unlearning):** The explorer model suppresses the spurious answer but preserves the true answer:

$$\begin{aligned} \pi_{\text{explorer}}(y_{\text{sp}}) &\ll \pi_{\text{anchor}}(y_{\text{sp}}), \\ \frac{\pi_{\text{explorer}}(y_{\text{true}})}{\pi_{\text{anchor}}(y_{\text{true}})} &\gg \frac{\pi_{\text{explorer}}(y_{\text{sp}})}{\pi_{\text{anchor}}(y_{\text{sp}})} \end{aligned}$$

The ratio inequality implies the true answer is more robust to unlearning.

**Assumption 3 (Large-Sample Consistency):** For sufficiently large  $G$ , by the Law of Large Numbers:

$$p_0(a) \xrightarrow{G \rightarrow \infty} \pi_{\text{anchor}}(a), \quad p_1(a) \xrightarrow{G \rightarrow \infty} \pi_{\text{explorer}}(a)$$

### D.3 Main Result

**Theorem:** Under Assumptions 1-3, DCRL selects the true answer ( $y_{\text{DC}}^* = y_{\text{true}}$ ), while majority vote selects the spurious answer ( $\hat{y}_{\text{MV}} = y_{\text{sp}}$ ). Thus,  $\text{Acc}(y_{\text{DC}}^*) > \text{Acc}(\hat{y}_{\text{MV}})$ .

### D.4 Proof

**Part 1: Majority Vote Converges to Spurious Answer :** By Assumption 1,  $\pi_{\text{anchor}}(y_{\text{sp}}) \gg \pi_{\text{anchor}}(y_{\text{true}})$ . By Assumption 3,  $p_0(y_{\text{sp}}) > p_0(a), \forall a \in \mathcal{A}$ . Thus  $\hat{y}_{\text{MV}} = \arg \max_a p_0(a) = y_{\text{sp}}$ .

**Part 2: Dual Consensus Converges to True Answer :** We show  $S(y_{\text{true}}) > S(y_{\text{sp}})$ . For large  $G$ :

$$S(a) \rightarrow \tilde{S}(a) = \frac{2\pi_{\text{anchor}}(a)\pi_{\text{explorer}}(a)}{\pi_{\text{anchor}}(a) + \pi_{\text{explorer}}(a)}.$$

Define robustness ratios (Assumption 2,  $r_{\text{true}} \gg r_{\text{sp}}$ ):

$$r_{\text{sp}} = \frac{\pi_{\text{explorer}}(y_{\text{sp}})}{\pi_{\text{anchor}}(y_{\text{sp}})}, \quad r_{\text{true}} = \frac{\pi_{\text{explorer}}(y_{\text{true}})}{\pi_{\text{anchor}}(y_{\text{true}})},$$

where  $r_{\text{true}} \gg r_{\text{sp}} \rightarrow 0$ .

Substitute ratios into  $\tilde{S}(a)$ :

$$\begin{aligned} \tilde{S}(y_{\text{sp}}) &= \frac{2\pi_{\text{anchor}}(y_{\text{sp}}) \cdot r_{\text{sp}}\pi_{\text{anchor}}(y_{\text{sp}})}{\pi_{\text{anchor}}(y_{\text{sp}}) + r_{\text{sp}}\pi_{\text{anchor}}(y_{\text{sp}})} \\ &= \frac{2r_{\text{sp}} \cdot \pi_{\text{anchor}}(y_{\text{sp}})}{1 + r_{\text{sp}}}. \end{aligned}$$

$$\begin{aligned} \tilde{S}(y_{\text{true}}) &= \frac{2\pi_{\text{anchor}}(y_{\text{true}}) \cdot r_{\text{true}}\pi_{\text{anchor}}(y_{\text{true}})}{\pi_{\text{anchor}}(y_{\text{true}}) + r_{\text{true}}\pi_{\text{anchor}}(y_{\text{true}})} \\ &= \frac{2r_{\text{true}} \cdot \pi_{\text{anchor}}(y_{\text{true}})}{1 + r_{\text{true}}}. \end{aligned}$$

By Assumption 2,  $r_{\text{sp}} \rightarrow 0 \implies \tilde{S}(y_{\text{sp}}) \rightarrow 0$ . Since  $\pi_{\text{anchor}}(y_{\text{true}}) > 0$  and  $r_{\text{true}} > 0$ , we have  $\tilde{S}(y_{\text{true}}) > 0$ . Thus  $\tilde{S}(y_{\text{true}}) > \tilde{S}(y_{\text{sp}})$ , so  $y_{\text{DC}}^* = \arg \max_a S(a) = y_{\text{true}}$ .

Dual Consensus enforces a robustness constraint—valid answers must be supported by both anchor and explorer, eliminating spurious answers fragile to unlearning and outperforming Majority Vote.

## E The Use of Large Language Models

We used large language models (LLMs) only to polish writing and improve textual clarity. No LLM was applied to research idea generation, experimental design, data analysis, or result derivation. All scientific contributions—conceptualization, methodology, experiments, and conclusions—were independently developed by the authors in full.