

# Understanding and Preventing Entropy Collapse in RLVR with On-Policy Entropy Flow Optimization

Anonymous ACL submission

## Abstract

Reinforcement learning with verifiable rewards (RLVR) has become an effective paradigm for improving the reasoning ability of large language models. However, widely used RLVR algorithms, such as GRPO, often suffer from entropy collapse, leading to premature determinism and unstable optimization. Existing remedies, including entropy regularization and ratio-based clipping heuristics, either control entropy in a coarse-grained manner or rely on approximate on-policy training. In this paper, we revisit entropy collapse from a token-level entropy flow perspective. Our analysis reveals that entropy-decreasing tokens consistently outweigh entropy-increasing ones, resulting in a severely imbalanced entropy flow. This perspective provides a unified explanation of entropy collapse in existing RLVR algorithms and highlights the importance of balancing entropy dynamics. Motivated by this analysis, we propose On-Policy Entropy Flow Optimization (OPEFO), an adaptive entropy flow balancing mechanism that rescales entropy-increasing and entropy-decreasing updates according to their contributions to entropy change, while remaining strict on-policy. Experiments on six mathematical reasoning benchmarks demonstrate that OPEFO improves training stability and final performance. We will release the code and models upon publication.

## 1 Introduction

Reinforcement Learning with Verifiable Rewards (RLVR) has emerged as an effective paradigm to advance reasoning capabilities in Large Language Models (LLMs, Lambert et al., 2024; Jaech et al., 2024; Guo et al., 2025; Yang et al., 2025a; Team et al., 2025). RLVR optimizes LLMs outputs via RL objectives guided by automated verifiable reward signals. However, recent RLVR algorithms, such as Proximal Policy Optimization (PPO, Schulman et al., 2017) and Group Relative Policy Optimization (GRPO, Shao et al., 2024), often suffer

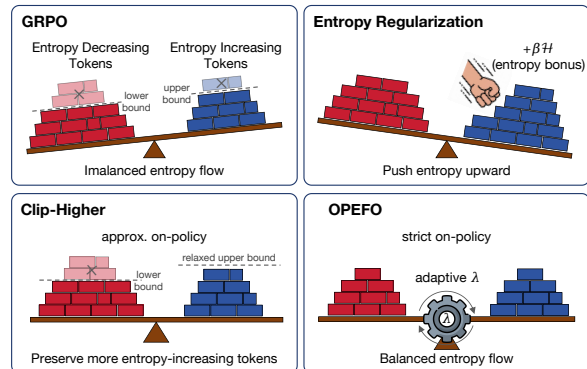


Figure 1: Entropy control mechanisms from the entropy flow perspective. GRPO suffers from imbalanced entropy flow. Entropy regularization increases entropy by adding an explicit entropy bonus; Clip-higher preserves more entropy-increasing tokens via relaxing the upper clipping bound. Differently, our OPEFO adaptively balances entropy flow via an adaptive scaling mechanism.

from *Entropy Collapse*: policy entropy drops in the early stage of training and continues to decline towards near zero. This indicates the policy becomes prematurely deterministic, limiting policy exploration and thus hindering LLMs’ reasoning capabilities (Cui et al., 2025; Hao et al., 2025b).

Recent approaches to mitigating entropy collapse fall into two categories, each with notable limitations, as shown in Figure 1. The first adopts entropy regularization, which introduces an explicit entropy bonus to encourage higher policy entropy (Mnih et al., 2016; Haarnoja et al., 2018). But this approach indiscriminately increases entropy and may cause excessive entropy growth in later training stages (Shen, 2025; Cheng et al., 2025). Besides, the entropy term is optimized independently of the policy-gradient objective, so it may dominate the advantage signal, leading to unstable updates and degraded performance (Zhang et al., 2025; Liu et al., 2025). The second category employs ratio- or clipping-based heuristics (Yu et al., 2025; Yang et al., 2025b). For instance, Clip-

066 higher (Yu et al., 2025) relaxes the upper clipping  
067 bound of the importance sampling ratio to preserve  
068 more entropy-increasing tokens. But they are only  
069 approximately on-policy due to their reliance on  
070 an outdated reference policy, making them theoret-  
071 ically less grounded than strict on-policy optimiza-  
072 tion (Baird et al., 1995; Sutton et al., 1999; Hao  
073 et al., 2025a; Zheng et al., 2025).

074 In this paper, we address these challenges from  
075 a novel perspective. We revisit entropy collapse  
076 through the lens of **entropy flow**: how entropy-  
077 increasing and entropy-decreasing token-level up-  
078 dates jointly shape overall entropy dynamics dur-  
079 ing training. Our analysis reveals that entropy-  
080 decreasing tokens dominate the early training  
081 phase, resulting in a severely imbalanced entropy  
082 flow with a strongly negative net entropy change.  
083 This imbalance provides a unified explanation for  
084 entropy collapse in existing RLVR algorithms.

085 To mitigate this imbalance issue, we propose **On-**  
086 **Policy Entropy Flow Optimization (OPEFO)**.  
087 OPEFO introduces an adaptive entropy flow balanc-  
088 ing mechanism that rescales the entropy-increasing  
089 and entropy-decreasing updates based on their re-  
090 spective contributions to entropy change. As such,  
091 our OPEFO adaptively balances the entropy flow  
092 and thus stabilizes policy entropy, as illustrated  
093 in Figure 1. Moreover, it remains strict on-policy  
094 since it avoids relying on reference policies as pre-  
095 vious methods. Empirical results demonstrate that  
096 OPEFO achieves best performance across two base  
097 models and six challenging mathematical reason-  
098 ing benchmarks. More importantly, it maintains the  
099 most stable entropy dynamics among strong base-  
100 lines. Overall, we summarize our contributions as  
101 follows:

- 102 • From an entropy flow perspective, we conduct  
103 a token-level analysis of entropy dynamics and  
104 show that entropy collapse can be interpreted  
105 as an imbalance between entropy-increasing and  
106 entropy-decreasing updates.
- 107 • We propose OPEFO, a strict on-policy entropy  
108 flow balancing mechanism that rescales entropy-  
109 increasing and entropy-decreasing token-level  
110 updates to stabilize policy entropy.
- 111 • We demonstrate through extensive experiments  
112 on mathematical reasoning benchmarks that  
113 OPEFO consistently improves training stabil-  
114 ity and final performance compared to existing  
115 RLVR methods.

## 2 Related Work 116

117 Reinforcement Learning with Verifiable Rewards  
118 (RLVR, Lambert et al., 2024) has recently emerged  
119 as an effective fine-tuning paradigm for large  
120 language models, achieving notable success in  
121 reasoning-intensive domains such as mathematics  
122 and programming (Shao et al., 2024; Guo et al.,  
123 2025; Yang et al., 2025a). Its core idea is to replace  
124 human feedback with automatically verifiable, ob-  
125 jective criteria as reinforcement learning rewards,  
126 thereby avoiding the cost and complexity of train-  
127 ing human preference models. OpenAI o1 (Ope-  
128 nAI, 2024) is among the first large-scale deploy-  
129 ments of this paradigm, demonstrating substantial  
130 performance gains on mathematical competition  
131 problems and code generation benchmarks. Fol-  
132 lowing this line of work, several subsequent mod-  
133 els—including DeepSeek-R1 (Guo et al., 2025),  
134 Kimi-1.5 (Team et al., 2025), and Qwen-2.5 (Yang  
135 et al., 2024)—have reported matched or improved  
136 results on similar reasoning benchmarks. Overall,  
137 RLVR has shown clear advantages over prior ap-  
138 proaches based solely on supervised fine-tuning for  
139 reasoning tasks.

140 These methods typically adopt GRPO-style train-  
141 ing, yet empirical studies consistently report the  
142 emergence of entropy collapse during optimiza-  
143 tion, where the policy’s entropy rapidly diminishes  
144 as training progresses. Early approaches draw on  
145 classical entropy regularization, introducing en-  
146 tropy bonuses or KL penalties to stabilize opti-  
147 mization (Mnih et al., 2016; Haarnoja et al., 2018).  
148 However, recent studies show that such techniques  
149 are highly sensitive to coefficient tuning in LLM  
150 and may even mislead optimization at critical states,  
151 yielding only coarse-grained effects on entropy con-  
152 trol (Cui et al., 2025; Shen, 2025).

153 More recent efforts attempt to mitigate entropy  
154 collapse by modifying key components of pol-  
155 icy optimization, including asymmetric or adap-  
156 tive ratio clipping (Yu et al., 2025; Yang et al.,  
157 2025b), selective optimization over high-entropy  
158 tokens (Wang et al., 2025b), or balancing pos-  
159 itive and negative samples (Zhu et al., 2025).  
160 Other methods introduce entropy-aware advan-  
161 tages or auxiliary objectives to encourage explo-  
162 ration (Cheng et al., 2025; Tan et al., 2025; Wang  
163 et al., 2025b,a; Deng et al., 2025). In this work, we  
164 adopt a token-level perspective to analyze entropy  
165 change and entropy flow, and study how these dy-  
166 namics evolve under strict on-policy optimization.

### 3 Preliminaries

In this section, we introduce the preliminaries of RLVR and policy entropy of LLMs.

#### 3.1 RLVR Algorithms

We consider reinforcement learning from verifiable rewards, where a policy model  $\pi_\theta$  autoregressively generates a token sequence  $y$  given a prompt  $x$ . The objective is to maximize the expected reward received from a verifier:

$$\mathcal{J}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} [r(x, y)] \quad (1)$$

where  $\mathcal{D}$  denotes the training data. Following the policy gradient theorem (Williams, 1992), the gradient of this objective can be estimated as:

$$\nabla_\theta \mathcal{J}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} \left[ \sum_{t=1}^{|y|} \nabla_\theta \log \pi_\theta(y_t | x, y_{<t}) \cdot A_t \right] \quad (2)$$

where  $A_t$  denotes the estimated advantage of token  $y_t$ . To avoid training an additional value network, recent algorithms such as GRPO (Shao et al., 2024) estimate token-level advantages via group-wise normalization:

$$A_t = \frac{r(y) - \text{mean}(r(y^{1:K}))}{\text{std}(r(y^{1:K}))} \quad (3)$$

where  $r(y^{1:K})$  denotes the rewards of  $K$  rollouts sampled for the same prompt. All tokens within the same response share the same normalized advantage value.

#### 3.2 Policy Entropy of LLMs

For an autoregressive language model  $\pi_\theta$ , uncertainty at each decoding step is characterized by the entropy of the next-token distribution. Given a state  $s_t = (x, y_{<t})$ , where  $a$  denotes a candidate next token sampled from the vocabulary, the token-level entropy  $\mathcal{H}_t$  is defined as:

$$\mathcal{H}_t = -\mathbb{E}_{a \sim \pi_\theta(\cdot|s_t)} [\log \pi_\theta(a|s_t)] \quad (4)$$

A smaller  $\mathcal{H}_t$  indicates higher confidence, while a larger value reflects greater uncertainty or exploration.

To capture uncertainty over an entire response and dataset  $\mathcal{D}$ , the policy entropy is defined as the average token entropy:

$$\mathcal{H}(\pi_\theta, \mathcal{D}) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)} \left[ \frac{1}{|y|} \sum_{t=1}^{|y|} \mathcal{H}_t \right] \quad (5)$$

While  $\mathcal{H}_t$  captures local uncertainty at individual decoding steps,  $\mathcal{H}(\pi_\theta, \mathcal{D})$  provides a global summary of the policy’s exploration level.

#### 3.3 Token-Level Entropy Change

While token entropy measures the model’s uncertainty at each decoding step, entropy collapse is driven by how entropy changes during policy updates. Following prior works (Hao et al., 2025b; Cui et al., 2025), we therefore analyze training dynamics through token-level entropy change, decomposing the overall entropy evolution into local changes at individual decoding steps. Directly computing the exact entropy change for large autoregressive models is intractable, due to complex dependencies across tokens. As a result, existing works commonly adopt a simplified tabular-softmax assumption, where each token’s logit is treated as conditionally independent.

**Assumption 1** (Parameter-independent softmax). Assume the policy  $\pi_\theta$  is a tabular softmax policy, where each state-action pair  $(s, a)$  is associated with an individual logit parameter  $z_{s,a}(\theta) = \theta_{s,a}$ .

**Theorem 1** (First-order entropy change). Under Assumption 1, the change of conditional entropy between two update steps is defined as  $\Delta \mathcal{H}_t \triangleq \mathcal{H}(\pi_\theta^{k+1} | s_t) - \mathcal{H}(\pi_\theta^k | s_t)$ . Then the first-order estimation of  $\Delta \mathcal{H}_t$  is:

$$\Delta \mathcal{H}_t = -\eta \mathbb{E}_{a \sim \pi_\theta^k(\cdot|s_t)} \left[ A_t (1 - \pi_\theta^k(a | s_t))^2 (\log \pi_\theta^k(a | s_t) + \mathcal{H}(\pi_\theta^k | s_t)) \right] \quad (6)$$

where  $\eta$  is the learning rate, and  $k$  indexes the policy update step.

This expression is derived from the first-order entropy change analysis introduced in Hao et al. (2025b), utilizing a Taylor expansion of the conditional entropy around the current policy logits. Compared to prior formulations based on importance ratios, we present the expression under the strict on-policy setting, where expectations are taken with respect to the current policy  $\pi_\theta^k$  without reference policies or importance ratios.

We adopt  $\Delta \mathcal{H}_t$  as a diagnostic quantity rather than a learning objective. Prior work shows that this first-order approximation closely tracks the exact entropy change in practice, justifying its use for analyzing training dynamics (Hao et al., 2025b). In our setting, it enables a token-level decomposition of entropy evolution, which we leverage to identify entropy imbalance and motivate our on-policy entropy flow balancing mechanism.

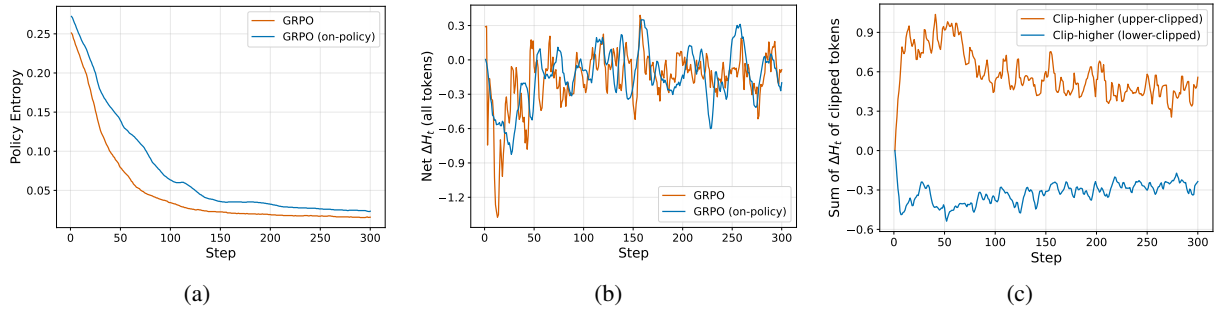


Figure 2: Empirical analysis of entropy dynamics in GRPO and Clip-higher. (a) Policy entropy over training steps for GRPO and its variant. (b) Net  $\Delta\mathcal{H}_t$  per update step for GRPO and its variant. (c) Sum of  $\Delta\mathcal{H}_t$  of upper- and lower-clipped tokens under Clip-higher.

## 4 Empirical Observations

In this section, we investigate how entropy evolves during RLVR training from an entropy flow perspective. We examine GRPO and its strict on-policy variant. GRPO performs 8 updates per batch, whereas the strict on-policy setting performs only one update per batch. In addition, we analyze Clip-higher as a representative clipping-based heuristic to understand, through token-level entropy change, why such methods can alleviate entropy collapse. All experiments are conducted using Qwen-2.5-Math-7B (Yang et al., 2024) as the base model, trained on DAPO-17k (Yu et al., 2025).

**Entropy collapse under GRPO.** Figure 2 (a) plots the evolution of policy entropy for GRPO and its strict on-policy variant. Both curves exhibit highly similar dynamics: starting at an entropy of approximately 0.25, dropping sharply within the first 50 training steps, and rapidly converging to nearly zero. This observation shows that entropy collapse is a common phenomenon in GRPO training, and motivates the need for a mechanism that can stabilize entropy dynamics.

**Token-level entropy change in GRPO.** We compute the aggregate token-level entropy change  $\sum_t \Delta\mathcal{H}_t$  at each training step, based on the Eq. 6. Figure 2 (b) shows that this quantity is strongly negative during the first 50 steps, indicating that dynamics consistently bias entropy downward. This behavior directly explains the sharp entropy drop observed in Figure 2(a). As training proceeds, the net entropy change approaches zero, coinciding with the stabilization of policy entropy. Together, these observations indicate that the evolution of policy entropy is driven by the cumulative token-level entropy changes, highlighting the importance of stabilizing this quantity throughout training.

**Clipping behavior of Clip-higher.** Since Clip-higher has demonstrated strong empirical effectiveness in alleviating entropy collapse, we analyze it as a representative clipping-based heuristic to understand its behavior from a token-level entropy flow perspective. For each training step, we compute the total  $\Delta\mathcal{H}_t$  contributed by tokens clipped at the upper and lower bounds, respectively. Figure 2 (c) reveals a clear pattern: tokens clipped at the upper bound tend to exhibit positive entropy change ( $\Delta\mathcal{H}_t > 0$ ), whereas those clipped at the lower bound typically exhibit negative entropy change ( $\Delta\mathcal{H}_t < 0$ ). This pattern provides a mechanistic explanation for the effectiveness of Clip-higher: by relaxing the upper bound, it selectively preserves entropy-increasing updates, thereby partially counteracting entropy collapse. However, as this effect relies on importance-ratio clipping, it is fundamentally incompatible with strict on-policy training.

**Summary.** These observations suggest that entropy collapse in RLVR can be interpreted from an entropy flow perspective, where entropy-decreasing updates consistently outweigh entropy-increasing ones, leading to a strongly negative net entropy change. Since clipping-based heuristics such as Clip-higher rely on importance-ratio clipping, they are incompatible with strict on-policy training, motivating the need for a direct mechanism to stabilize entropy dynamics under strict on-policy optimization.

## 5 On-Policy Entropy Flow Optimization

To address the imbalanced entropy flow observed in Section 4, we propose On-Policy Entropy Flow Optimization (OPEFO), a mechanism designed for stabilizing entropy dynamics under strict on-policy RLVR training.

## 5.1 Entropy Flow Decomposition

To formalize the entropy flow, we decompose entropy evolution into token-level contributions within each policy update. For each generated token  $y_t$ , we estimate the corresponding entropy change  $\Delta\mathcal{H}_t$  (Eq. 6) induced by a policy-gradient update. Using the sign of  $\Delta\mathcal{H}_t$ , we partition the token-level updates within a batch into two sets:

- **Entropy-increasing set ( $\mathcal{S}^+$ ):** Tokens with  $\Delta\mathcal{H}_t > 0$ , corresponding to updates that broaden the model’s predictive distribution.
- **Entropy-decreasing set ( $\mathcal{S}^-$ ):** Tokens with  $\Delta\mathcal{H}_t < 0$ , corresponding to updates that sharpen the distribution.

As observed in Section 4, GRPO-style training exhibits a pronounced imbalance in entropy flow, where entropy-decreasing updates consistently outweigh entropy-increasing ones across training steps,  $\sum_{t \in \mathcal{S}^-} |\Delta\mathcal{H}_t| > \sum_{t \in \mathcal{S}^+} \Delta\mathcal{H}_t$ . As this imbalance persists across updates, the policy entropy progressively collapses.

## 5.2 Balanced On-Policy Objective

To stabilize entropy dynamics under strict on-policy RLVR training, we introduce a reweighted on-policy gradient objective that rescales updates from the two token sets via a balancing coefficient  $\lambda$ :

$$\nabla_{\theta} J_{\text{OPEFO}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(\cdot|x)} \left[ (1 + \lambda) \sum_{t \in \mathcal{S}^+} \nabla_{\theta} \log \pi_{\theta}(y_t|x, y_{<t}) \cdot A_t + (1 - \lambda) \sum_{t \in \mathcal{S}^-} \nabla_{\theta} \log \pi_{\theta}(y_t|x, y_{<t}) \cdot A_t \right] \quad (7)$$

where  $\lambda \in (-1, 1)$  controls the balance between entropy-increasing and entropy-decreasing updates.

This objective offers two practical advantages. First, it is fully compatible with strict on-policy training, as it does not rely on importance sampling ratios or reference policies. Second, by only rescaling gradient components, it preserves the original optimization direction without introducing auxiliary entropy bonuses.

## 5.3 Adaptive Entropy Flow Scaling

The remaining question is how to determine an appropriate value of  $\lambda$ , which we compute adap-

tively from batch-level entropy statistics. Specifically, we adopt zero entropy flow as a local stabilizing criterion at the batch level. This criterion is motivated by the observation that entropy collapse arises when entropy-decreasing updates consistently outweigh entropy-increasing ones. By enforcing a balance between these two opposing forces, we establish a condition for stabilizing entropy dynamics without external intervention. Formally, we seek a value of  $\lambda$  such that:

$$\sum_{t \in \mathcal{S}^+} (1 + \lambda) \Delta\mathcal{H}_t + \sum_{t \in \mathcal{S}^-} (1 - \lambda) \Delta\mathcal{H}_t \approx 0 \quad (8)$$

Solving this constraint yields a unique solution:

$$\lambda^* = \frac{\sum_{t \in \mathcal{S}^-} |\Delta\mathcal{H}_t| - \sum_{t \in \mathcal{S}^+} \Delta\mathcal{H}_t}{\sum_{t \in \mathcal{S}^-} |\Delta\mathcal{H}_t| + \sum_{t \in \mathcal{S}^+} \Delta\mathcal{H}_t} \quad (9)$$

This closed-form expression provides an adaptive entropy flow balancing coefficient computed directly from the current batch. We do not claim  $\lambda^*$  to be globally optimal across tasks or training stages; rather, it acts as a self-adjusting mechanism that compensates for transient entropy imbalances.

## 5.4 Practical Implementation

OPEFO is straightforward to implement and requires only minimal modification to standard on-policy RLVR pipelines. Practically, implementing OPEFO requires only a few lines of modification to standard GRPO code: compute  $\Delta\mathcal{H}_t$ , group tokens by sign, compute  $\lambda^*$  using Eq. 9, and reweight gradients accordingly. This procedure is fully compatible with strict on-policy training and introduces negligible computational overhead. Figure 3 shows the implementation code of OPEFO.

```
def compute_opefo_loss(log_prob,
    advantages, delta_H, eps=1e-12, **
    args):
    pg_terms = - log_prob * advantages
    S_pos = delta_H > 0
    S_neg = delta_H < 0

    pos_mag = delta_H[S_pos].sum()
    neg_mag = delta_H[S_neg].abs().sum()

    denom = (neg_mag + pos_mag).
    clamp_min(eps)
    lambda_s = (neg_mag - pos_mag) /
    denom

    pg_terms[S_pos] *= (1.0 + lambda_s)
    pg_terms[S_neg] *= (1.0 - lambda_s)
    return pg_terms.mean()
```

Figure 3: Implementation code of OPEFO loss.

## 6 Experiments

### 6.1 Experimental setup

**Training.** Following prior work (Hao et al., 2025b; Cui et al., 2025), we use Qwen2.5-Math-7B (Yang et al., 2024) as the primary base model, and additionally include Qwen3-4B-Base (Yang et al., 2025a) to evaluate the generality of our method. The training codebase is adapted from Verl (Sheng et al., 2025). The training data is DAPO-17K (Yu et al., 2025), which contains high-quality reasoning trajectories annotated with verifiable rewards.

For rollout, the prompt batch size is set to 32, with a single gradient update performed per rollout step under the strict on-policy setting. We use the AdamW (Loshchilov and Hutter, 2017) optimizer with a constant learning rate of  $2.83e-6$ <sup>1</sup>, together with a linear warm-up over the first 10 rollout steps. The maximum response length is set to 3000 tokens for Qwen2.5-Math-7B-base and 8000 tokens for Qwen3-4B-base. All experiments are conducted on 32 A100 GPUs.

**Baselines.** We compare OPEFO with several representative baselines. GRPO is used as the primary reference method, with both the upper and lower clipping bounds set to 0.2. Entropy regularization (Entropy-Reg) augments the objective with an explicit entropy regularization term, where the coefficient is fixed to 0.01. Clip-higher relaxes the upper clipping bound to 0.28 while keeping the lower bound unchanged at 0.2, following the standard configuration in prior work. In addition, Clip-Cov and KL-Cov are included, using the settings reported in Cui et al. (2025). Except for GRPO (strict on-policy), all baseline methods are trained using 8 gradient updates per batch with a learning rate of  $1e-6$ .

**Evaluation.** We evaluate models on six widely used mathematical reasoning benchmarks: AIME24, AIME25, AMC23 (Li et al., 2024), MATH500 (Hendrycks et al., 2021), Minerva Math (Lewkowycz et al., 2022), and Olympiad-Bench (He et al., 2024). Following prior work (Hao et al., 2025b), validation is performed with a top-p value of 0.7 and temperature 1.0 across all models and test sets. We report Avg@32 for AIME24 (30

<sup>1</sup>Following common practice, the learning rate is scaled proportionally to the square root of the effective batch size increase (i.e.,  $\sqrt{8}$ ) to maintain a comparable optimization noise scale (Goyal et al., 2017; Smith and Le, 2017).

problems), AIME25 (30 problems), and AMC23 (40 problems) due to their smaller size, and report Avg@1 for all other benchmarks. All evaluations are zero-shot with no additional prompts. All methods save a checkpoint every 10 steps, and the checkpoint achieving the highest average accuracy (Avg.) across benchmarks is selected for evaluation.

### 6.2 Main results

Table 1 reports the results of two base models on six mathematical reasoning benchmarks. We highlight three key observations.

First, among a set of competitive baselines, OPEFO achieves the best overall performance. On Qwen2.5-Math-7B and Qwen3-Base-4B, OPEFO attains average accuracies of 52.4% and 51.9%, respectively, outperforming the second-best methods by margins of 1.7% and 1.8%. These gains are observed across models and datasets, indicating that OPEFO improves reasoning ability in a robust manner.

Second, comparing GRPO with its strict on-policy variant reveals a clear and consistent empirical advantage of the strict on-policy setting over the approximate on-policy one. This provides empirical support for the claim discussed earlier that avoiding stale reference policies and distribution shifts can improve policy optimization in RLVR. However, strict on-policy GRPO alone, despite achieving higher accuracy, does not prevent entropy collapse, as shown in subsequent analyses.

Third, under the same strict on-policy setting, OPEFO further improves performance over strict on-policy GRPO by explicitly balancing entropy-increasing and entropy-decreasing updates, yielding average gains of 2.3% and 1.8% on Qwen2.5-Math-7B and Qwen3-Base-4B, respectively. This highlights the complementary role of entropy flow balancing: while strict on-policy training provides a cleaner optimization signal, OPEFO addresses the remaining instability observed under strict on-policy GRPO, and we observe additional performance improvements under this stabilized entropy dynamics.

### 6.3 Training Dynamics Analysis

In this section, we conduct an empirical analysis of the training dynamics of Qwen2.5-Math-7B under several representative training methods, examining how these methods affect entropy, training reward, and response length over training.

Method	AIME24	AIME25	AMC23	MATH500	Minerva	Olympiad	Avg.
<i>Qwen2.5-Math-7B</i>	13.8	5.3	44.6	39.6	9.9	13.8	21.2
GRPO	26.7	13.0	69.8	80.3	37.1	46.1	45.5
GRPO (Strict on-policy)	32.6	15.4	81.8	83.5	38.9	48.1	50.1
Entropy-Reg	31.7	13.3	74.4	81.9	40.5	45.1	47.8
Clip-higher	30.5	17.4	78.5	83.5	39.7	49.4	49.8
Clip-Cov	32.2	18.5	77.9	85.1	40.3	50.2	50.7
KL-Cov	32.6	17.9	78.3	84.6	40.9	48.7	50.5
<b>OPEFO</b>	<b>34.5</b>	<b>19.2</b>	<b>82.2</b>	<b>85.3</b>	<b>41.6</b>	<b>51.8</b>	<b>52.4</b>
<i>Qwen3-Base-4B</i>	9.7	8.8	51.1	75.4	33.1	40.5	36.4
GRPO	19.4	17.3	66.1	80.6	36.3	49.4	44.8
GRPO (Strict on-policy)	26.2	22.3	71.5	86.3	38.9	55.7	50.1
Entropy-Reg	22.7	20.8	70.9	83.4	37.1	53.3	48.0
Clip-higher	23.5	21.7	70.3	84.8	39.2	55.3	49.1
Clip-Cov	24.2	23.1	71.9	85.7	39.7	56.9	48.9
KL-Cov	24.7	22.4	72.3	86.1	40.3	55.8	49.2
<b>OPEFO</b>	<b>27.7</b>	<b>24.8</b>	<b>73.1</b>	<b>87.5</b>	<b>41.2</b>	<b>57.4</b>	<b>51.9</b>

Table 1: Main results on mathematical reasoning benchmarks. All results are presented as percentages. The best are in **bold**.

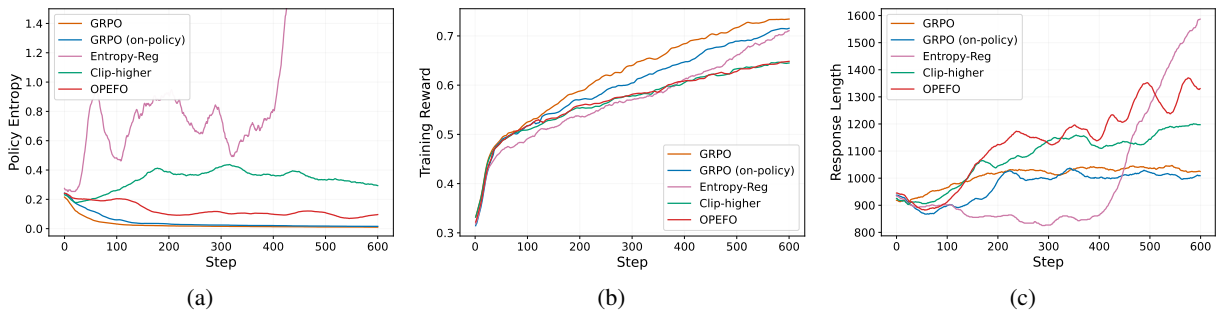


Figure 4: Training dynamics under different methods: (a) policy entropy, (b) training reward, and (c) response length.

**Overall Entropy.** As shown on the Figure 4 (a), both standard GRPO and its strict on-policy variant exhibit a rapid collapse of policy entropy during training. While strict on-policy GRPO achieves higher accuracy than the standard setting (Table 1), it does not inherently prevent entropy collapse. In contrast, entropy regularization leads to uncontrolled entropy growth in later training stages (after approximately 400 steps), whereas Clip-higher partially mitigates collapse but still suffers from noticeable oscillations. By comparison, OPEFO maintains a smooth and stable entropy trajectory, avoiding both premature determinism and uncontrolled entropy inflation.

**Training Reward.** Figure 4 (b) shows steady upward trends across all methods. Among them, GRPO and its strict on-policy variant achieve the highest training reward while also exhibiting the

lowest policy entropy. When interpreted together with the entropy curves, this pattern is consistent with premature exploitation of a limited set of high-reward answer patterns, rather than sustained exploration throughout training. A more desirable objective is to improve training performance while maintaining sufficient entropy to preserve diversity in the learned policy. OPEFO follows this direction by stabilizing entropy while still achieving competitive reward growth.

**Response Length.** Figure 4 (c) shows the evolution of response length during training. For GRPO and its on-policy variant, the response length quickly saturates and stops increasing after around 200 steps, while Clip-higher stabilizes at a later stage, around 300 steps. Entropy regularization exhibits a sharp increase in response length after around 400 steps, coinciding with the surge

Method	AIME24	AIME25	AMC23	MATH500	Minerva	Olympiad	Avg.
GRPO (Strict on-policy)	26.7	13.0	69.8	80.3	37.1	46.1	45.5
Static Scaling	30.0	16.5	81.0	83.8	38.0	49.5	49.8
One-side ( $\mathcal{S}^+$ only)	32.8	18.2	83.1	84.2	38.7	47.5	50.8
One-side ( $\mathcal{S}^-$ only)	32.3	16.3	83.6	83.9	38.3	47.2	50.3
<b>OPEFO</b>	<b>34.5</b>	<b>19.2</b>	<b>82.2</b>	<b>85.3</b>	<b>41.6</b>	<b>51.8</b>	<b>52.4</b>

Table 2: Ablation of balancing coefficient  $\lambda^*$  on Qwen2.5-Math-7B. GRPO (Strict on-policy) corresponds to  $\lambda^* = 0$ . “Static Scaling” uses a fixed  $\lambda^* = 0.001$ . “One-side ( $\mathcal{S}^+$  only)” and “One-side ( $\mathcal{S}^-$  only)” apply  $\lambda^*$  to entropy-increasing and entropy-decreasing updates, respectively. OPEFO applies  $\lambda^*$  to both sides.

in policy entropy observed in Figure 4 (a), indicating that the model begins to generate longer but increasingly unconstrained responses. In contrast, OPEFO continues to produce longer and controlled responses throughout training. We do not claim longer responses are inherently better; rather, under the verifiable-reward setting, response length serves as a behavioral indicator of whether the model continues to explore more intermediate reasoning paths.

#### 6.4 Analysis of the Balancing Coefficient $\lambda^*$

We analyze the role of the balancing coefficient  $\lambda^*$  on Qwen2.5-Math-7B, examining both different balancing strategies and the behavior of the analytically derived  $\lambda^*$  during training.

**Ablation on  $\lambda^*$ .** OPEFO adaptively reweights updates by amplifying entropy-increasing updates in  $\mathcal{S}^+$  and attenuating entropy-decreasing updates in  $\mathcal{S}^-$  using a dynamically computed  $\lambda^*$  (Eq. 9). We first compare OPEFO against a static scaling baseline with a fixed  $\lambda = 0.001$ , chosen as the average value of  $\lambda^*$  observed over the course of training. As shown in Table 2, static scaling yields moderate improvements over strict on-policy GRPO ( $\lambda^* = 0$ ), but consistently underperforms OPEFO, highlighting the importance of dynamically adjusting  $\lambda^*$ .

We further consider two one-sided variants that apply  $\lambda^*$  only to entropy-increasing updates ( $\mathcal{S}^+$  only) or only to entropy-decreasing updates ( $\mathcal{S}^-$  only). Both variants outperform strict on-policy GRPO, but remain consistently inferior to full OPEFO, suggesting that jointly balancing leads to better overall performance.

**Dynamics of  $\lambda^*$ .** Beyond ablations on balancing strategies, we further examine how the balancing coefficient  $\lambda^*$  evolves during training. As shown in Figure 5,  $\lambda^*$  exhibits a non-stationary pattern over training. Overall,  $\lambda^*$  remains positive throughout training, indicating a persistent need to encourage

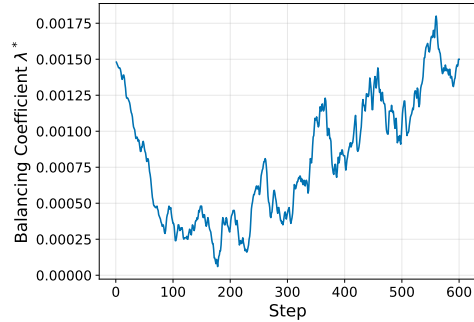


Figure 5: Evolution of the balancing coefficient  $\lambda^*$  over training steps.

entropy-increasing updates to counteract entropy collapse. We observe that  $\lambda^*$  first decreases in the early stage and then gradually increases later in training, reflecting different entropy balancing demands across training phases. Taken together,  $\lambda^*$  displays bounded oscillations rather than converging to a fixed value, reflecting an adaptive balancing mechanism that adjusts to non-stationary entropy flow during training.

## 7 Conclusion

We proposed On-Policy Entropy Flow Optimization (OPEFO), a strict on-policy mechanism for stabilizing entropy dynamics in RLVR training. By analyzing entropy collapse from a token-level entropy flow perspective, we showed that it arises from a persistent imbalance between entropy-increasing and entropy-decreasing updates. OPEFO directly addresses this issue by balancing the two opposing entropy flows without relying on reference policies or heuristic entropy regularization. Experiments on six mathematical reasoning benchmarks across two base models demonstrate that OPEFO consistently improves training stability and final performance over strong RLVR baselines. These results suggest entropy flow balancing as a simple and effective principle for stabilizing reasoning-oriented reinforcement learning.

## 596 Limitations

597 Our method offers a principled and interpretable en-  
598 tropy flow perspective for stabilizing RLVR train-  
599 ing under a strict on-policy setting, and demon-  
600 strates strong empirical performance across reason-  
601 ing benchmarks. We nevertheless note several  
602 limitations.

603 First, our analysis relies on a first-order approxi-  
604 mation of token-level entropy change under a sim-  
605 plified softmax assumption, which captures domi-  
606 nant entropy trends but abstracts away higher-order  
607 interactions in large transformer models.

608 Second, the proposed entropy flow balancing  
609 mechanism operates as a local, batch-level stabi-  
610 lizer, where the zero entropy flow criterion serves  
611 as a sufficient condition rather than a globally op-  
612 timal objective, and entropy flow behaviors may  
613 vary across tasks.

614 Finally, our evaluation focuses on strict on-  
615 policy RLVR for mathematical reasoning. While  
616 OPEFO is conceptually applicable to other do-  
617 mains and reward structures, a systematic study  
618 under dense rewards or more complex credit as-  
619 signment settings is left for future work.

620 We view these limitations not as fundamental  
621 constraints of the proposed method, but as oppor-  
622 tunities for extending and refining entropy-aware  
623 optimization methods in future RLVR systems.

## 624 Acknowledgments

625 We only use AI-assisted tools for language polish-  
626 ing and improving clarity.

## 627 References

628 Leemon Baird and 1 others. 1995. [Residual algorithms:](#)  
629 [Reinforcement learning with function approximation.](#)  
630 In *Proceedings of the twelfth international confer-*  
631 *ence on machine learning*, pages 30–37.

632 Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai,  
633 Wayne Xin Zhao, Zhenliang Zhang, and Furu Wei.  
634 2025. [Reasoning with exploration: An entropy per-](#)  
635 [spective.](#) *arXiv preprint arXiv:2506.14758*.

636 Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan  
637 Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen  
638 Fan, Huayu Chen, Weize Chen, and 1 others. 2025.  
639 [The entropy mechanism of reinforcement learning](#)  
640 [for reasoning language models.](#) *arXiv preprint*  
641 *arXiv:2505.22617*.

642 Jia Deng, Jie Chen, Zhipeng Chen, Wayne Xin Zhao,  
643 and Ji-Rong Wen. 2025. [Decomposing the entropy-](#)  
644 [performance exchange: The missing keys to unlock-](#)

[ing effective reinforcement learning.](#) *arXiv preprint*  
*arXiv:2508.02260*. 645 646

Priya Goyal, Piotr Dollár, Ross Girshick, Pieter No-  
ordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew  
Tulloch, Yangqing Jia, and Kaiming He. 2017. [Ac-](#)  
647 [curate, large minibatch sgd: Training imagenet in 1](#)  
648 [hour.](#) 649 650 651

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*. 652 653 654 655 656 657

Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and  
Sergey Levine. 2018. [Soft actor-critic: Off-policy](#)  
[maximum entropy deep reinforcement learning with](#)  
[a stochastic actor.](#) In *International conference on*  
*machine learning*, pages 1861–1870. Pmlr. 658 659 660 661 662

Yaru Hao, Li Dong, Xun Wu, Shaohan Huang, Zewen  
Chi, and Furu Wei. 2025a. [On-policy rl with optimal](#)  
[reward baseline.](#) *arXiv preprint arXiv:2505.23585*. 663 664 665

Zhezhen Hao, Hong Wang, Haoyang Liu, Jian Luo,  
Jiarui Yu, Hande Dong, Qiang Lin, Can Wang, and  
Jiawei Chen. 2025b. [Rethinking entropy interven-](#)  
[tions in rlvr: An entropy change perspective.](#) *arXiv*  
*preprint arXiv:2510.10150*. 666 667 668 669 670

Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding  
Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han,  
Yujie Huang, Yuxiang Zhang, and 1 others. 2024.  
[Olympiadbench: A challenging benchmark for pro-](#)  
[moting agi with olympiad-level bilingual multimodal](#)  
[scientific problems.](#) In *Proceedings of the 62nd An-*  
*ual Meeting of the Association for Computational*  
*Linguistics (Volume 1: Long Papers)*, pages 3828–  
3850. 671 672 673 674 675 676 677 678 679

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul  
Arora, Steven Basart, Eric Tang, Dawn Song, and Ja-  
cob Steinhardt. 2021. [Measuring mathematical prob-](#)  
[lem solving with the math dataset.](#) *arXiv preprint*  
*arXiv:2103.03874*. 680 681 682 683 684

Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richard-  
son, Ahmed El-Kishky, Aiden Low, Alec Helyar,  
Aleksander Madry, Alex Beutel, Alex Carney, and 1  
others. 2024. [Openai o1 system card.](#) *arXiv preprint*  
*arXiv:2412.16720*. 685 686 687 688 689

Nathan Lambert, Jacob Morrison, Valentina Pyatkin,  
Shengyi Huang, Hamish Ivison, Faeze Brahman,  
Lester James V Miranda, Alisa Liu, Nouha Dziri,  
Shane Lyu, and 1 others. 2024. [Tulu 3: Pushing fron-](#)  
[tiers in open language model post-training.](#) *arXiv*  
*preprint arXiv:2411.15124*. 690 691 692 693 694 695

Aitor Lewkowycz, Anders Andreassen, David Dohan,  
Ethan Dyer, Henryk Michalewski, Vinay Ramasesh,  
Ambrose Slone, Cem Anil, Imanol Schlag, Theo  
Gutman-Solo, and 1 others. 2022. [Solving quan-](#)  
[titative reasoning problems with language models.](#) 696 697 698 699 700

701	<i>Advances in neural information processing systems</i> ,	Hongze Tan, Jianfei Pan, Jinghao Lin, Tao Chen, Zhi-	753
702	35:3843–3857.	hang Zheng, Zhihao Tang, and Haihua Yang. 2025.	754
703	Jia Li, Edward Beeching, Lewis Tunstall, Ben Lip-	<a href="#">Gtpo and grpo-s: Token and sequence-level re-</a>	755
704	kin, Roman Soletskyi, Shengyi Huang, Kashif Rasul,	<a href="#">ward shaping with policy entropy.</a>	756
705	Longhui Yu, Albert Q Jiang, Ziju Shen, and 1 oth-	<i>arXiv preprint</i>	757
706	ers. 2024. <a href="#">Numinamath: The largest public dataset</a>	<i>arXiv:2508.04349.</i>	
707	<a href="#">in ai4maths with 860k pairs of competition math</a>	Kimi Team, Angang Du, Bofei Gao, Bowei Xing,	758
708	<a href="#">problems and solutions.</a>	Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun	759
709	<i>Hugging Face repository</i> ,	Xiao, Chenzhuang Du, Chonghua Liao, and 1 others.	760
	13(9):9.	2025. <a href="#">Kimi k1. 5: Scaling reinforcement learning</a>	761
710	Jiacai Liu. 2025. <a href="#">How does rl policy entropy converge</a>	<a href="#">with llms.</a>	762
711	<a href="#">during iteration.</a>	<i>arXiv preprint arXiv:2501.12599.</i>	
	<i>Zhihu Zhuanlan.</i>	Jiakang Wang, Runze Liu, Fuzheng Zhang, Xiu Li, and	763
712	Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin	Guorui Zhou. 2025a. <a href="#">Stabilizing knowledge, promot-</a>	764
713	Dong, Yejin Choi, Jan Kautz, and Yi Dong. 2025.	<a href="#">ing reasoning: Dual-token constraints for rlvr.</a>	765
714	<a href="#">Prorl: Prolonged reinforcement learning expands rea-</a>	<i>arXiv</i>	766
715	<a href="#">soning boundaries in large language models.</a>	<i>preprint arXiv:2507.15778.</i>	
716	<i>arXiv preprint arXiv:2505.24864.</i>	Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shix-	767
717	Ilya Loshchilov and Frank Hutter. 2017. <a href="#">Decou-</a>	uan Liu, Rui Lu, Kai Dang, Xionghui Chen, Jianxin	768
718	<a href="#">pled weight decay regularization.</a>	Yang, Zhenru Zhang, and 1 others. 2025b. <a href="#">Beyond</a>	769
719	<i>arXiv preprint</i>	<a href="#">the 80/20 rule: High-entropy minority tokens drive</a>	770
	<i>arXiv:1711.05101.</i>	<a href="#">effective reinforcement learning for llm reasoning.</a>	771
		<i>arXiv preprint arXiv:2506.01939.</i>	772
720	Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi	Ronald J Williams. 1992. <a href="#">Simple statistical gradient-</a>	773
721	Mirza, Alex Graves, Timothy Lillicrap, Tim Harley,	<a href="#">following algorithms for connectionist reinforcement</a>	774
722	David Silver, and Koray Kavukcuoglu. 2016. <a href="#">Asyn-</a>	<a href="#">learning.</a>	775
723	<a href="#">chronous methods for deep reinforcement learning.</a>	<i>Machine learning</i> , 8(3):229–256.	
724	In <i>International conference on machine learning</i> ,	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	776
725	pages 1928–1937. PmLR.	Binyuan Hui, Bo Zheng, Bowen Yu, Chang	777
726	OpenAI. 2024. <a href="#">Learning to reason with llms.</a>	Gao, Chengen Huang, Chenxu Lv, and 1 others.	778
		2025a. <a href="#">Qwen3 technical report.</a>	779
		<i>arXiv preprint</i>	780
		<i>arXiv:2505.09388.</i>	
727	John Schulman, Filip Wolski, Prafulla Dhariwal,	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui,	781
728	Alec Radford, and Oleg Klimov. 2017. <a href="#">Proxi-</a>	Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu,	782
729	<a href="#">mal policy optimization algorithms.</a>	Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jian-	783
730	<i>arXiv preprint</i>	hong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang,	784
	<i>arXiv:1707.06347.</i>	Jingren Zhou, Junyang Lin, Kai Dang, and 22 oth-	785
		ers. 2024. <a href="#">Qwen2.5 technical report.</a>	786
731	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	<i>arXiv preprint</i>	787
732	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	<i>arXiv:2412.15115.</i>	
733	Zhang, YK Li, Yang Wu, and 1 others. 2024.	Shihui Yang, Chengfeng Dou, Peidong Guo, Kai	788
734	<a href="#">Deepseekmath: Pushing the limits of mathematical</a>	Lu, Qiang Ju, Fei Deng, and Rihui Xin. 2025b.	789
735	<a href="#">reasoning in open language models.</a>	<a href="#">Dcpo: Dynamic clipping policy optimization.</a>	790
736	<i>arXiv preprint</i>	<i>arXiv preprint arXiv:2509.02333.</i>	791
	<i>arXiv:2402.03300.</i>	Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan,	792
737	Han Shen. 2025. <a href="#">On entropy control in llm-rl algo-</a>	Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan,	793
738	<a href="#">rithms.</a>	Gaohong Liu, Lingjun Liu, and 1 others. 2025. <a href="#">Dapo:</a>	794
	<i>arXiv preprint arXiv:2509.03493.</i>	<a href="#">An open-source llm reinforcement learning system</a>	795
739	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin	<a href="#">at scale.</a>	796
740	Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin	<i>arXiv preprint arXiv:2503.14476.</i>	
741	Lin, and Chuan Wu. 2025. <a href="#">Hybridflow: A flexible</a>	Ruipeng Zhang, Ya-Chien Chang, and Sicun Gao. 2025.	797
742	<a href="#">and efficient rlhf framework.</a>	<a href="#">When maximum entropy misleads policy optimiza-</a>	798
743	In <i>Proceedings of the Twentieth European Conference on Computer Sys-</i>	<a href="#">tion.</a>	799
744	<i>tems</i> , pages 1279–1297.	<i>arXiv preprint arXiv:2506.05615.</i>	
745	Samuel L Smith and Quoc V Le. 2017. <a href="#">A bayesian</a>	Haizhong Zheng, Jiawei Zhao, and Beidi Chen. 2025.	800
746	<a href="#">perspective on generalization and stochastic gradient</a>	<a href="#">Prosperity before collapse: How far can off-policy</a>	801
747	<a href="#">descent.</a>	<a href="#">rl reach with stale data on llms?</a>	802
	<i>arXiv preprint arXiv:1710.06451.</i>	<i>arXiv preprint</i>	803
		<i>arXiv:2510.01161.</i>	
748	Richard S Sutton, David McAllester, Satinder Singh,	Xinyu Zhu, Mengzhou Xia, Zhepei Wei, Wei-Lin Chen,	804
749	and Yishay Mansour. 1999. <a href="#">Policy gradient methods</a>	Danqi Chen, and Yu Meng. 2025. <a href="#">The surprising</a>	805
750	<a href="#">for reinforcement learning with function approxima-</a>	<a href="#">effectiveness of negative reinforcement in llm reason-</a>	806
751	<a href="#">tion.</a>	<a href="#">ing.</a>	807
752	<i>Advances in neural information processing</i>	<i>arXiv preprint arXiv:2506.01347.</i>	
	<i>systems</i> , 12.		

## A Theorem Proof Details 808

The following derivation is adapted from prior work (Hao et al., 2025b) and is included for reference and completeness. 809  
810

**Theorem 1** (First-order entropy change). Under Assumption 1, the change of conditional entropy between two update steps is defined as  $\Delta \mathcal{H}_t \triangleq \mathcal{H}(\pi_\theta^{k+1} | s_t) - \mathcal{H}(\pi_\theta^k | s_t)$ . Then the first-order estimation of  $\Delta \mathcal{H}_t$  is: 811  
812  
813

$$\Delta \mathcal{H}_t = -\eta \mathbb{E}_{a \sim \pi_\theta^k(\cdot | s_t)} [A_t (1 - \pi_\theta^k(a | s_t))^2 (\log \pi_\theta^k(a | s_t) + \mathcal{H}(\pi_\theta^k | s_t))] \quad (10) \quad 814$$

where  $\eta$  is the learning rate, and  $k$  indexes the policy update step. Note that compared to the formulation in Hao et al. (2025b), Eq. 10 is specialized to the strict on-policy setting: the weight term  $w_t$  (defined there as  $w_t = \mathbb{I}_{\text{clip}} r_t A_t$ ) reduces to the plain advantage  $A_t$  since no importance ratio or clipping is used in our updates. 815  
816  
817  
818

**Proof.** The proof is similar to that of (Liu, 2025). Taking the first-order Taylor expansion, we have 819

$$\Delta H_t \triangleq \mathcal{H}(\pi_\theta^{k+1} | s_t) - \mathcal{H}(\pi_\theta^k | s_t) \approx \langle \nabla_\theta \mathcal{H}(\pi_\theta^k | s_t), z^{k+1} - z^k \rangle \quad 820$$

Since we have the log trick  $\mathbb{E}_{a \sim \pi_\theta(\cdot | s)} [\nabla_\theta \log \pi_\theta(a | s)] = 0$ , the gradient term can be derived as 821

$$\begin{aligned} \nabla_\theta \mathcal{H}(\pi_\theta | s) &= \nabla_\theta \mathcal{H}(\pi_\theta(\cdot | s)) \\ &= \nabla_\theta (-\mathbb{E}_{a \sim \pi_\theta(\cdot | s)} [\log \pi_\theta(a | s)]) \\ &= -\mathbb{E}_{a \sim \pi_\theta(\cdot | s)} [\nabla_\theta \log \pi_\theta(a | s) + \log \pi_\theta(a | s) \nabla_\theta \log \pi_\theta(a | s)] \\ &= -\mathbb{E}_{a \sim \pi_\theta(\cdot | s)} [\log \pi_\theta(a | s) \nabla_\theta \log \pi_\theta(a | s)]. \end{aligned} \quad 822$$

Then we have 823

$$\begin{aligned} \Delta H_t &= \langle \nabla_\theta \mathcal{H}(\theta^k | s_t), (z^{k+1} - z^k) \rangle \\ &= -\left\langle \mathbb{E}_{a \sim \pi_\theta^k(\cdot | s_t)} [\log \pi_\theta(a | s_t) \nabla_\theta \log \pi_\theta(a | s_t)], \theta^{k+1} - \theta^k \right\rangle \\ &= -\mathbb{E}_{a \sim \pi_\theta^k(\cdot | s_t)} \left[ \log \pi_\theta(a | s_t) \langle \nabla_\theta \log \pi_\theta(a | s_t), \theta^{k+1} - \theta^k \rangle \right] \\ &= -\mathbb{E}_{a \sim \pi_\theta^k(\cdot | s_t)} \left[ \log \pi_\theta(a | s_t) \sum_{a' \in \mathcal{A}} \frac{\partial \log \pi_\theta(a | s_t)}{\partial \theta_{s_t, a'}} (\theta_{s_t, a'}^{k+1} - \theta_{s_t, a'}^k) \right] \\ &= -\mathbb{E}_{a \sim \pi_\theta^k(\cdot | s_t)} \left[ \log \pi_\theta(a | s_t) \sum_{a' \in \mathcal{A}} (\mathbf{1}\{a = a'\} - \pi(a' | s_t)) (\theta_{s_t, a'}^{k+1} - \theta_{s_t, a'}^k) \right] \\ &= -\mathbb{E}_{a \sim \pi_\theta^k(\cdot | s_t)} \left[ \left( \log \pi_\theta(a | s_t) - \mathbb{E}_{\tilde{a} \sim \pi_\theta^k(\cdot | s_t)} \log \pi_\theta(\tilde{a} | s_t) \right) \right. \\ &\quad \left. \left( \theta_{s_t, a}^{k+1} - \theta_{s_t, a}^k - \mathbb{E}_{a' \sim \pi_\theta^k(\cdot | s_t)} (\theta_{s_t, a'}^{k+1} - \theta_{s_t, a'}^k) \right) \right] \\ &= -\mathbb{E}_{a \sim \pi_\theta^k(\cdot | s)} \left[ [\log \pi_\theta^k(a | s) + \mathcal{H}(\cdot | s)] \left[ (1 - \pi_\theta^k(a | s)) (z_{s_t, a}^{k+1} - z_{s_t, a}^k) \right] \right] \\ &= -\mathbb{E}_{a \sim \pi_\theta^k(\cdot | s)} \left[ [\log \pi_\theta^k(a | s) + \mathcal{H}(\cdot | s)] \left[ w(s | a) (1 - \pi_\theta^k(a | s))^2 \right] \right], \end{aligned} \quad 824$$

where  $w(s | a)$  is the weight in the policy gradient. □ 825

**B Detailed Information For Test Dataset**

<b>Test Datasets</b>	<b>#Questions</b>	<b>Level</b>
AIME24	30	Olympiad
AIME25	30	Olympiad
AMC23	40	Intermediate
MATH500	500	Advanced
Minerva	272	Graduate
OlympiadBench	675	Olympiad

Table 3: Dataset statistics.