

# QaRL: Rollout-Aligned Quantization-Aware RL for Fast and Stable Training under Training–Inference Mismatch

Anonymous ACL submission

## Abstract

Large language model (LLM) reinforcement learning (RL) pipelines are often bottlenecked by rollout generation, making end-to-end training slow. Recent work mitigates this by running rollouts with quantization to accelerate decoding, which is the most expensive stage of the RL loop. However, these setups destabilize optimization by amplifying the training-inference gap: rollouts are operated at low precision, while learning updates are computed at full precision. To address this challenge, we propose QaRL (Rollout Alignment Quantization-Aware RL), which aligns training-side forward with the quantized rollout to minimize mismatch. We further identify a failure mode in quantized rollouts: long-form responses tend to produce repetitive, garbled tokens (error tokens). To mitigate these problems, we introduce TBPO (Trust-Band Policy Optimization), a sequence-level objective with dual clipping for negative samples, aimed to keep updates within the trust region. On Qwen3-30B-A3B MoE for math problems, QaRL outperforms quantized-rollout training by +5.5 while improving stability and preserving low-bit throughput benefits.

## 1 Introduction

Recent reasoning LLMs, such as OpenAI o1 (Jaeuch et al., 2024), have demonstrated the effectiveness of Chain-of-Thought (CoT) (Wei et al., 2022) for solving complex problems. More recently, DeepSeek-R1 (Guo et al., 2025) has shown that reinforcement learning with simple rule-based reward functions (RLVR), can induce emergent reasoning behaviors and yield gains in challenging domains such as math problem solving (Luo et al., 2025). A key driver behind these improvements is test-time scaling (Muennighoff et al., 2025): reasoning models often generate longer CoT, trade additional computation for accuracy. However, longer generations increase

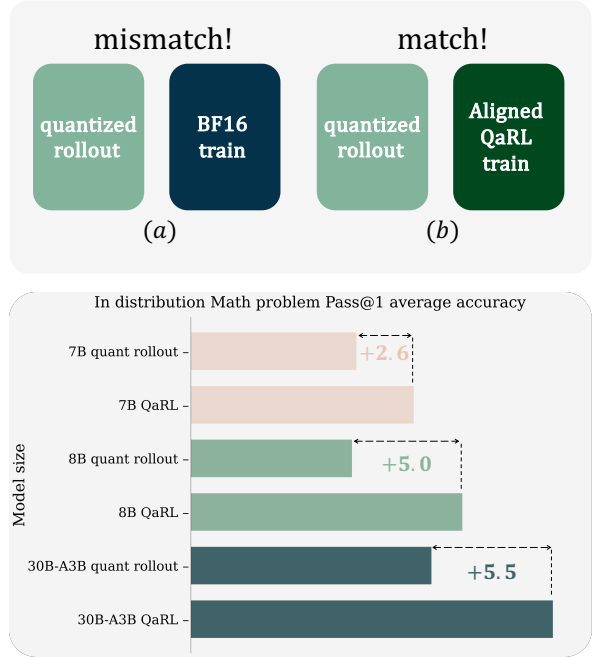


Figure 1: (a) Quantized rollout alone introduces mismatch. (b) Aligned QaRL alleviates it, improving performance across model sizes, including MoE.

training-time cost, since RL continuously samples long responses during optimization.

Unlike SFT, which requires only a single forward pass, a standard RL step (Sheng et al., 2024; Hu et al., 2025) for LLMs involves three phases: (i) rollout to generate responses, (ii) a forward pass to compute token probabilities, and (iii) a backward pass to update model via policy gradients. Autoregressive decoding generates response token by token, making rollouts the dominant cost, comprising roughly 70% of RL training time. Consequently, quantizing rollout model is a natural way to accelerate RL optimization (Liu et al., 2025b). However, it also creates a training-inference mismatch challenge: responses are sampled from low bit rollout, while updates are performed by full precision learner.

To address both efficiency and stability, in this work, we study **quantization-aware RL**. By exe-

cutting rollouts under low-bit quantization, we substantially reduce the dominant generation cost and accelerate end-to-end RL training. To mitigate the resulting training-inference mismatch, we further perform rollout-aligned quantization-aware training on the learner side, aligning learners policy with the quantized behavior used for sampling.

Moreover, we find a critical failure mode in quantized rollouts: noise accumulates over long generations, producing off-trajectory repetitive and garbled tokens. These **error tokens** are typically assigned very low probability under the policy, driving the policy ratio (and mismatch reweighting) to extreme. Such outliers are not reliably controlled by standard PPO-style clipping, breaking the intended trust region behavior and destabilizing training. To address this, we introduce a trust-band control: we apply dual clipping tailored to negative samples and perform sequence-level objectives on both the policy ratio and the mismatch weight, drop entire responses that exceed the bands.

Extensive experiments on math, logic and code benchmarks demonstrate that QaRL matches BF16 training performance and outperform quantized rollout training. Even on unstable MoE models such as Qwen3-30B-A3B-Base, QaRL achieves an average math score of 51.2, close to BF16’s 52.1 and surpassing quantized rollout’s 45.7, while delivering a  $1.3\times$  training speedup over BF16. Our contributions are as follows:

1. We present a practical quantized aware RL pipeline for decoupled, hybrid RL systems, and introduce **QaRL** to minimize the mismatch between quantized rollouts and learner-side training.
2. We identify error tokens as the key cause of training instability under quantized rollouts, and propose **Trust Band Policy Optimization (TBPO)**, a sequence-level dual-clipping strategy that keeps updates within a trust region and stabilizes convergence.

## 2 Preliminaries

### 2.1 Group Relative Policy Optimization.

GRPO (Shao et al., 2024) is a PPO style (Schulman et al., 2017) policy gradient objective that removes the value function (critic) and uses group relative rewards as the baseline. Given a query  $q$ , we sample a group of  $G$  responses  $\{o_1, \dots, o_G\}$  from the old policy  $o_i \sim \pi_{\theta_{\text{old}}}(\cdot | q)$ . Let  $r_i$  be the

scalar (verifiable) reward for  $o_i$  and  $|o_i|$  its length. GRPO defines a group-normalized advantage  $A_i$  (broadcast to tokens as  $A_{i,t} = A_i$ ) and the token-level importance ratio  $R_{i,t}(\theta)$ :

$$A_i = \frac{r_i - \text{mean}\{r_1, \dots, r_G\}}{\text{std}\{r_1, \dots, r_G\}},$$

$$R_{i,t}(\theta) = \frac{\pi_{\theta}(o_{i,t} | q_i)}{\pi_{\theta_{\text{old}}}(o_{i,t} | q_i)}$$

To prevent overly large updates and ensure the update remains within the trust region (Schulman et al., 2015), GRPO applies PPO-style clipping and defines the token-level clipped policy loss:

$$\mathcal{L}_{i,t}(\theta) = -\min\left(R_{i,t}(\theta) A_i, \tilde{R}_{i,t}(\theta) A_i\right),$$

where  $\tilde{R}_{i,t} = \text{clip}(R_{i,t}, 1 - \epsilon, 1 + \epsilon)$  and  $\epsilon$  controls the trust region. The GRPO optimization objective is defined as follows:

$$\mathcal{L}_{\text{GRPO}}(\theta) = \mathbb{E}_{q, o \sim \pi_{\theta_{\text{old}}}} \left[ \frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \mathcal{L}_{i,t}(\theta) \right].$$

### 2.2 Quantization

Quantization accelerates LLM inference by enabling low-bit GEMM operations, where weights and activations are compressed into low bit representations that modern GPUs can process directly via Tensor Core kernels.

**Integer quantization.** A full-precision tensor  $W \in \mathbb{R}$  is mapped to  $W_q \in \mathbb{Z}$  via scale  $s$  and zero-point  $z$ :

$$\text{Quant} : W_q = \text{clamp}\left(\left\lfloor \frac{W}{s} \right\rfloor + z, q_{\min}, q_{\max}\right),$$

$$\text{Dequant} : \hat{W} = s(W_q - z),$$

where  $\lfloor \cdot \rfloor$  denotes rounding. Quantized GEMM computes the integer matmul and rescales the result (accumulating in higher precision):

$$Y \approx s_x s_w \cdot ((X_q - z_x)(W_q - z_w)). \quad (1)$$

**Floating-point quantization.** For low-bit floating formats (e.g., FP8/FP4), values are cast to reduced precision with scaling:

$$\text{Quant} : s = \frac{\max(|W|)}{\alpha}, W_q = \lfloor W/s \rfloor_{\text{FP}k},$$

$$\text{Dequant} : \hat{W} = s \cdot W_q,$$

where  $\lfloor \cdot \rfloor_{\text{FP}k}$  denotes casting to FP $k$  and  $\alpha$  is the maximum finite representable value ( $\alpha = 6$

for e2m1FP4,  $\alpha = 448$  for e4m3FP8). Low-bit GEMM performs multiplication on quantized operands and rescales:

$$Y \approx (s_x s_w) \cdot (X_q W_q), \quad (2)$$

with  $X_q, W_q$  as the low-bit floating representations. We denote  $x$ -bit weight and  $y$ -bit activation quantization as  $W \times A_y$ .

**Quantization-Aware Training (QAT).** To recover accuracy lost in low-bit conversion, QAT integrates quantization noise into the training loop (Liu et al., 2024). Specifically, it injects fake quant into the forward pass to simulate rounding, while master weights remain in full precision. The Straight-Through Estimator (STE) approximates gradients for non-differentiable backward pass.

**Fully Quantized Training (FQT).** Distinct from QAT, which employs fake quant to merely simulate quantization noise while maintaining high-precision arithmetic, Fully Quantized Training (or Low-Bit Training) fundamentally executes operations using actual low-bit data types (e.g., FP8, INT8) (Xi et al., 2024; Zhao et al., 2025). By leveraging low-precision hardware kernels for both forward and backward passes, this paradigm directly reduces the computational overhead and memory footprint of the training process itself.

### 2.3 Training–Inference Mismatch in RL

Modern LLM RL pipelines are typically hybrid systems: for throughput, rollouts are generated by a high-performance inference engine (e.g., vLLM, SGLang), while policy optimization are carried out by a training backend (e.g., FSDP/Megatron). Although both components load the same parameter  $\theta_{\text{old}}$ , the rollout and the training engine may produce different results for the same query ( $\text{LLM}_{\text{rollout}}(q) \neq \text{LLM}_{\text{train}}(q)$ ). This occurs because they often implement different kernels (e.g., different attention kernels, batch variant operations). As a result, the rollout engine induces a **sampler policy**  $\pi_{\text{sampler}}(\cdot | \theta_{\text{old}})$  that is not exactly the same distribution as the **learner policy**  $\pi_{\text{learner}}(\cdot | \theta_{\text{old}})$ . We refer to this discrepancy as **training–inference (learner–sampler) mismatch**.

Consider PPO’s clipped objective:

$$\mathbb{E}_{a \sim \pi_{\theta_{\text{old}}}} \left[ \min \left( r_{\text{prox}} \hat{A}, \text{clip}(r_{\text{prox}}, 1 - \epsilon, 1 + \epsilon) \hat{A} \right) \right]$$

where proximal importance sampling  $r_{\text{prox}} = \frac{\pi(a|\theta)}{\pi(a|\theta_{\text{old}})}$  and  $\hat{A}$  is an advantage estimate. In a hybrid system, however, tokens are sampled from

$\pi_{\text{sampler}}(\cdot | \theta_{\text{old}})$ , while  $r_{\text{prox}}$  is computed using learner’s distributions:

$$\mathbb{E}_{a \sim \pi_{\text{sampler}}(\theta_{\text{old}})} \left[ \min \left( \frac{\pi_{\text{learner}}(a | \theta)}{\pi_{\text{learner}}(a | \theta_{\text{old}})} \hat{A}, \text{clip} \left( \frac{\pi_{\text{learner}}(a | \theta)}{\pi_{\text{learner}}(a | \theta_{\text{old}})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A} \right) \right]$$

Consequently, the trust region is enforced on the sampler distribution rather than the learner distribution, causing the update to deviate from its intended target. Drawing inspiration from Decoupled PPO (Hilton et al., 2022), where a mismatch between the behavior and proximal policies is corrected via importance sampling, we can reweight the updates by multiplying with the ratio:

$$w_{\text{mismatch}} = \frac{\pi_{\text{learner}}(a | \theta_{\text{old}})}{\pi_{\text{sampler}}(a | \theta_{\text{old}})} \quad (3)$$

Applying  $w_{\text{mismatch}}$  yields

$$\mathbb{E}_{a \sim \pi_{\text{sampler}}(\theta_{\text{old}})} \left[ \frac{\pi_{\text{learner}}(a | \theta_{\text{old}})}{\pi_{\text{sampler}}(a | \theta_{\text{old}})} \cdot \min \left( r_{\text{prox}} \hat{A}, \text{clip}(r_{\text{prox}}, 1 - \epsilon, 1 + \epsilon) \hat{A} \right) \right]. \quad (4)$$

Intuitively, this correction is conceptually orthogonal to PPO proximal ratio  $r_{\text{prox}} = \frac{\pi_{\text{learner}}(a|\theta)}{\pi_{\text{learner}}(a|\theta_{\text{old}})}$ , while  $w_{\text{mismatch}}$  rectifies the distribution shift caused by the system-level training–inference mismatch. Together, they ensure policy optimization remains within the learner’s trust region. Recent works have further explored this direction, such as **TIS** (Yao et al., 2025) (which truncates upper bound of  $w_{\text{mismatch}}$ ) and **MIS** (Liu et al., 2025a) (reject samples with overly large  $w_{\text{mismatch}}$  and apply sequence level  $w_{\text{mismatch}}$ ).

## 3 Methods

### 3.1 Rollout-Aligned Quantization Aware Reinforcement Learning

To alleviate the decoding bottleneck during rollouts, a practical approach is to quantize the rollout model. By lowering the precision of weights and activations (e.g., W8A8 or W4A16), we can significantly accelerate rollout generation. However, quantizing only the rollout engine (we term this **Quantized Rollout Training**) introduces a more severe form of training–inference mismatch.

In this regime, the sampler policy  $\pi_{\text{sampler}}$  becomes a quantized approximation  $\pi_{\text{quant-sampler}}$ , while the learner policy  $\pi_{\text{learner}}$  typically remains

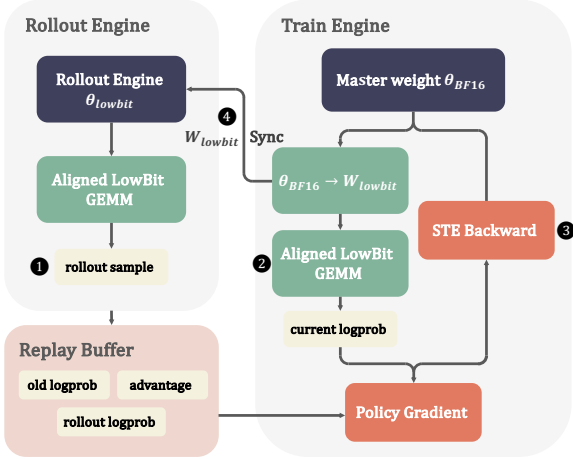


Figure 2: **Overview of the QaRL pipeline in a hybrid RL system.** ① The quantized rollout engine  $\theta_{lowbit}$  generates samples. ② The training engine maintains  $\theta_{BF16}$  master weights and performs rollout-aligned low-bit GEMM to compute current logprob. ③ Policy gradients are computed using replay buffer data to update the model via STE. ④ The updated low-bit weights  $W_{lowbit}$  are synchronized to the rollout engine.

full-precision  $\pi_{BF16\text{-learner}}$ . The resulting mismatch makes the importance weight

$$w_{\text{mismatch}} = \frac{\pi_{\text{learner}}}{\pi_{\text{sampler}}} = \frac{\pi_{BF16\text{-learner}}}{\pi_{\text{quant-sampler}}}$$

drift far away from 1.00 as Fig. 4(d). When both sampler and learner operate in full precision, the mismatch is typically limited to implementation details (e.g. different kernels), so the distribution shift is mild and can often be compensated by reweighting with  $w_{\text{mismatch}}$ . In contrast, quantized rollout training introduces larger distribution gap and  $w_{\text{mismatch}}$  become more extreme. This effect is especially pronounced for long responses, quantization-induced errors accumulate across decoding steps, making the divergence grow with response length and increasingly difficult to correct. To fundamentally resolve this mismatch, we align the learner with the quantized sampler by adopting **Quantization-Aware Reinforcement Learning (QaRL)**.

But building a quantized training pipeline in a decoupled hybrid RL stack is non-trivial: rollouts are generated by an inference engine, while policy optimization runs in a separate training backend. To address this challenge, we decompose our implementation into three components: (1) **low-bit rollout inference**, (2) **rollout aligned quantization aware training**, and (3) **weight synchronization from training to inference**. Concretely:

### 1. Quantized inference in the rollout engine.

We deploy a low-bit rollout model in the inference engine to accelerate decoding. Empirically, for larger models (especially MoE), we observe that W4A16 can deliver higher throughput than W8A8 due to memory constraints. Therefore, we support both configurations to maximize rollout efficiency across different model regimes.

### 2. Rollout aligned quantization aware training.

On the training side, we maintain master weights  $\theta_{BF16}$  in high-precision. Standard QAT simulates quantization error by inserting fake-quant operators  $[\text{dequant}[\text{quant}(\theta_{BF16})]]$ , while computing in high precision. To minimize the gap between  $\pi_{\text{quant-sampler}}$  and  $\pi_{\text{quant-learner}}$ , we perform on-the-fly quantization during the forward pass:

$$\hat{W}_{\text{low-bit}} = \text{Quant}(\theta_{BF16})$$

Critically, rather than merely simulating quantization, we execute low-bit GEMM directly on these quantized tensors, thereby precisely mirroring the arithmetic behavior of the rollout engine. During the backward pass, gradients are computed via the Straight-Through Estimator (STE) and applied to the master weights  $\theta_{BF16}$ . Our approach lies between fake quant QAT and end-to-end low-precision FQT: we execute the forward pass in low bit while keeping the backward pass in full precision. We avoid FQT for larger quantization error introduced by 4-bit gradients.

### 3. Weight updates from training to inference.

After each optimization step, the learners parameters change, requiring the rollout engine to be refreshed accordingly. Since the training engine already materializes  $\hat{W}_{\text{low-bit}}$  for the aligned low-bit GEMM, we directly publish these low-bit tensors to the inference engine at each step. This avoids redundant re-quantization and ensures the sampler weights remain in the exact low-bit format learner uses during forward. The overall pipeline is illustrated in Fig. 2.

Notably, unlike prior work (Wasti et al., 2025; SGLang Team, 2025) that enforces strictly bitwise consistent on-policy execution, our approach tolerates discrepancy between the sampler and learner policies  $\pi_{\text{quant-sampler}} \neq \pi_{\text{quant-learner}}$  and compensates for it using  $w_{\text{mismatch}}$ . Achieving bitwise alignment requires both batch and tensor-parallel invariance kernel, which is not a "free lunch", incurring an average  $2\times$  slower (He and Lab, 2025; Zhang et al., 2025). Consequently, rather

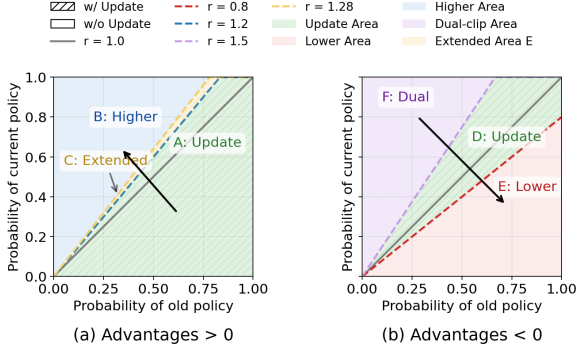


Figure 3: **Token level policy clipping regions.** Axes represent token probabilities under the old and current policies, with the slope defining the  $r_{\text{prox}} = \text{prob}_{\text{current}}/\text{prob}_{\text{old}}$ , with arrows indicating the direction of the policy update.

than bitwise-identical kernels, we focus on using aligned low-bit forward to ensure that optimization remains robustly within the intended trust region. As for master weight precision, previous work (Qi et al., 2025) claims that using the finer granularity of  $\theta_{\text{FP16}}$  can mitigate mismatch. However, we empirically find that using  $\theta_{\text{FP16}}$  with dynamic loss scaling causes gradient NaN underflow, and this phenomenon does not occur with  $\theta_{\text{BF16}}$ .

### 3.2 Trust Band Policy Optimization

Although rollout-aligned quantized training substantially mitigates instability arising from training–inference mismatch, we observe a remaining failure mode: under-trained quantized policy tend to produce **repetitive and garbled tokens** in long responses, which we term **error tokens**, illustrated in appendix table 6. This stems from the error-amplification dynamics of autoregressive decoding: an error token at step  $t$  sends the model off-trajectory, causing subsequent tokens generated from a corrupted state and amplifying degradation.

This phenomenon interacts closely with the trust-region control in PPO-style objectives. Standard PPO clipping operates directionally: for positive advantages, it constrains only the upper bound to  $(0, 1 + \epsilon)$ , while for negative advantages, it constrains only the lower bound to  $(1 - \epsilon, +\infty)$ . These prevent the updated policy drift too far from the old policy in the update direction. Recently, several works have revisited this design. For example, DAPO (Yu et al., 2025) proposes clipping higher for  $A > 0$  as C:extended region in Fig. 3(a), arguing that high-entropy which correspond to low-probability tokens, can encourage exploration.

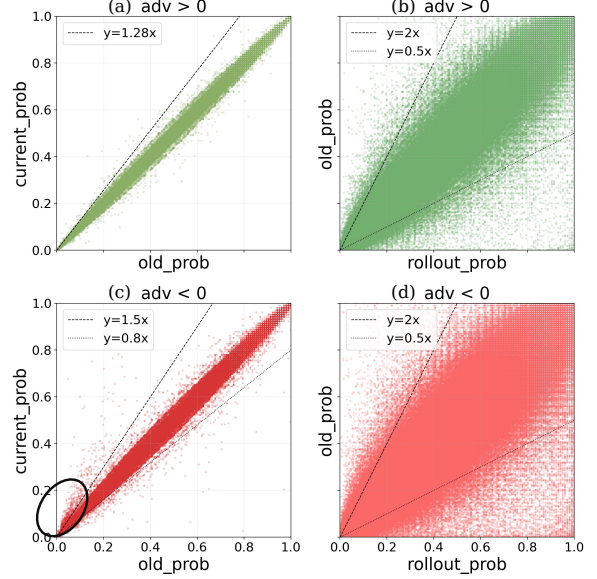


Figure 4: **Distribution of token level  $r_{\text{prox}}$  and  $w_{\text{mismatch}}$ .** Lines indicate clipping boundaries.

**Dual Clipping.** In our setting, the more critical issue emerges with negative samples. Empirically, error tokens are rare in positive trajectories but predominate in negative ones. As shown in Fig. 4(c), these negative error tokens are concentrated in the low-probability regions of both the old and current policies ( $\pi_{\text{old},t}$  and  $\pi_{\theta,t}$ ). In this regime, the PPO policy ratio becomes more sensitive to baseline probabilities:

$$r_{\text{prox},t} = \begin{cases} 0.80/0.77 \approx 1.03 & p_{\text{old}} = 0.77, p_{\theta} = 0.80 \\ 0.05/0.02 = 2.5 & p_{\text{old}} = 0.02, p_{\theta} = 0.05 \end{cases} \quad 341$$

Here, identical absolute perturbations produce dramatically different relative changes depending on  $\pi_{\text{old},t}$ . For typical high-confidence tokens (e.g.,  $p_{\text{old}} = 0.77 \rightarrow p_{\theta} = 0.80$ ), the ratio remains modest at  $\approx 1.07$ . Conversely, for low-probability tail tokens (e.g.,  $p_{\text{old}} = 0.02 \rightarrow p_{\theta} = 0.05$ ), the ratio explodes to 2.5. And low-prob tokens are inherently more prone to probability shifts. This effect is directly visible in Fig. 5(b): when  $\pi_{\text{old}}$  is small, negative-advantage tokens exhibit a much heavier tail of  $r_{\text{prox}}$ , producing extreme ratios. Under  $A_t < 0$ , standard PPO clipping is one-sided (it only enforces a lower bound), so these high  $r_{\text{prox}}$  tail tokens can dominate the gradient magnitude, inflate variance, and break the intended trust region. To mitigate this failure mode, we leverage a **negative dual clipping scheme** (Ye et al., 2020): for  $A < 0$ , we not only enforce the lower bound but also impose an explicit upper bound as

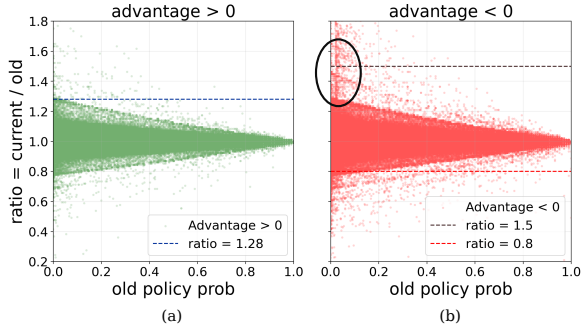


Figure 5: **Relationship between  $r_{\text{prox}}$  and old policy probability.** Low-probability negative tokens are more prone to extreme  $r_{\text{prox}}$ .

361 F: Dual region in Fig. 3(b), preventing exploding  
 362 ratios from low-probability error tokens from dom-  
 363 inating the update.

364 **Sequence Level Objectives.** Furthermore,  
 365 since RLVR optimize outcome rewards, GRPO’s  
 366 token-level importance sampling serves merely as  
 367 a first-order approximation. Token-level clipping  
 368 fails to rectify **mid-response errors**: once a single  
 369 token deviates, the subsequent tokens no longer lie  
 370 on the intended trajectory (Fig 6).

371 To remedy this, we return to a stricter trust-  
 372 region view and treat each entire response as single  
 373 action. Concretely, we apply sequence-level  
 374 ratios and clipping (Zheng et al., 2025), mask  
 375 out the entire responses from policy update when  
 376 their sequence-level ratio exceeds the bounds. An-  
 377 other advantage of sequence-level ratio is its abil-  
 378 ity to discriminate error tokens and low-prob  
 379 exploration tokens, which token-level bounds can-  
 380 not resolve. Error tokens typically occur in clus-  
 381 ters, causing high variance in  $r_{\text{seq-prox}}$  that readily  
 382 exceeds trust-region thresholds. In contrast, ex-  
 383 ploration tokens generate more stable sequence-  
 384 level ratios, as the geometric averaging over the  
 385 sequence reduces sensitivity to individual pertur-  
 386 bations.

387 Formally, we treat the entire response  $o =$   
 388  $(o_1, \dots, o_L)$  as a single action and defining its  
 389 sequence-level probability as the geometric mean  
 390 of the token probabilities. Consequently, we formu-  
 391 late  $r_{\text{seq-prox}}$  and  $w_{\text{seq-mismatch}}$ :

392 
$$\pi(o | q) \triangleq \exp\left(\frac{1}{L} \sum_{t=1}^L \log \pi(o_t | q, o_{<t})\right).$$

393 
$$r_{\text{seq-prox}}(\theta) \triangleq \frac{\pi_{\text{current learner}}(o | q)}{\pi_{\text{old learner}}(o | q)},$$

394 
$$w_{\text{seq-mismatch}}(\theta) \triangleq \frac{\pi_{\text{old learner}}(o | q)}{\pi_{\text{old sampler}}(o | q)}.$$



Figure 6: **A mid response error propagates to future tokens.** Although the initial garbled tokens are clipped, the repetitive tokens induced by this error are not clipped by token-level objectives.

For mismatch weight, we truncate with a two-  
 sided cap to control variance:

395 
$$\tilde{w}(\theta) \triangleq \text{clip}(w(\theta), [-\log c, \log c]), \quad (5)$$
 396 397

398 where  $c > 1$  is the TIS cap. For the proximal ratio,  
 399 positive-advantage samples use the standard PPO-  
 400 style bound  $[0, 1 + \epsilon_h]$ , while negative samples are  
 401 constrained by dual-sided band  $[1 - \delta_\ell, 1 + \delta_h]$ :

402 
$$\tilde{r}(\theta) = \begin{cases} \text{clip}(r(\theta), 0, 1 + \epsilon_h) & \text{if } \hat{A} \geq 0 \\ \text{clip}(r(\theta), 1 - \delta_\ell, 1 + \delta_h) & \text{if } \hat{A} < 0 \end{cases}$$
 403

404 Put everything together, the sequence-level surro-  
 405 gate (shared by all tokens in  $o$ ) is

406 
$$\mathcal{J}(\theta) = \mathbb{E}_{q, o \sim \pi_{\text{old sampler}}(\cdot | q)} \left[ \tilde{w}(\theta) \cdot \tilde{r}(\theta) \cdot \hat{A} \right].$$
 407

## 4 Experiments

408 **Training Setup.** We utilize OpenR1-Math-  
 409 46K (Yan et al., 2025), a dataset of 46,000  
 410 mathematical problems, and conduct experiments  
 411 on Qwen2.5-Math-1.5B/7B, Qwen3-8B-Base,  
 412 and Qwen3-30B-A3B-Base (Yang et al., 2025).  
 413 We employ GRPO with TIS for the BF16 RL  
 414 baseline and quantized training, while GSPO is  
 415 adopted for MoE models to enhance stability. The  
 416 Muon optimizer (Jordan et al., 2024) exhibits  
 417 significantly faster convergence than AdamW and  
 418 is therefore used across all experiments; detailed  
 419 hyperparameters are provided in Appendix B.

420 **Evaluation.** We evaluate on math benchmarks  
 421 (AIME2024/2025, AMC, Math-500, Olympiad-  
 422 Bench (He et al., 2024), Minerva) and out-of-  
 423 distribution benchmarks (ARC-Challenge (Clark  
 424 et al., 2018), GPQA-Diamond (Rein et al., 2024),  
 425 LiveCodeBench, MMLU Pro).

### 4.1 Main Results

426 Table 1 and Fig. 7 demonstrate a consistent trend  
 427 across model scales: quantized rollout training un-  
 428 dermines stability and yields lower final accuracy  
 429

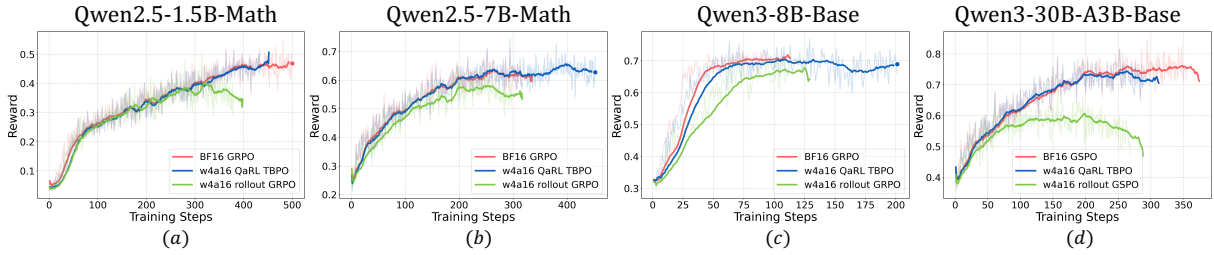


Figure 7: **Training reward curves across different models.** Our QaRL TBPO demonstrates stability over quantized rollout training, and converging to reward levels nearly identical to the full-precision BF16 baseline.

Model	In-Distribution Performance						Out-of-Distribution Performance					
	AIME 24/25	AMC	MATH-500	Minerva	Olympiad	Avg.	ARC-c	GPQA	MMLU-Pro	LiveCodeBench	Avg.	
Qwen2.5-1.5B-Math	4.5/2.8	26.5	50.8	21.6	20.3	21.0	11.7	12.4	10.4	2.7	9.3	
BF16 GRPO	12.5/9.2	43.5	71.4	36.2	34.5	34.5	58.9	27.8	26.7	11.9	31.3	
w4a16 rollout GRPO	7.9/6.4	38.1	65.3	30.2	28.9	29.3	50.1	19.8	21.2	4.3	23.8	
w4a16 QaRL TBPO	12.5/10.4	46.6	69.8	31.9	32.6	33.9	57.6	21.1	25.3	8.1	28.0	
Qwen2.5-7B-Math	15.0/6.4	46.2	67.4	32.9	23.9	33.4	62.6	28.5	32.1	8.0	32.8	
BF16 GRPO	19.5/11.6	59.6	80.0	45.9	43.4	43.3	80.4	38.2	46.3	14.6	44.8	
w4a16 rollout GRPO	19.1/9.6	54.6	79.1	42.8	40.6	40.9	73.9	31.6	40.0	8.7	38.5	
w4a16 QaRL TBPO	19.5/13.3	58.1	81.8	43.4	45.3	43.5	81.1	37.1	45.8	13.9	44.4	
Qwen3-8B-Base	7.9/9.6	46.3	74.2	42.7	39.3	36.6	44.9	31.2	49.5	23.0	37.1	
BF16 GRPO	28.3/19.3	64.3	88.1	54.5	56.7	51.8	93.0	46.3	65.1	45.7	62.5	
w4a16 rollout GRPO	20.0/12.5	53.0	82.2	50.7	45.1	43.9	91.6	43.9	61.5	41.4	59.5	
w4a16 QaRL TBPO	26.6/16.9	62.2	83.6	52.5	51.7	48.9	92.3	45.0	63.8	43.4	61.0	
Qwen3-30B-A3B-Base	15.4/7.9	49.0	67.4	31.2	38.1	34.8	61.3	34.8	52.5	28.5	44.2	
BF16 GSPO	27.9/21.6	63.2	88.8	54.7	56.7	52.1	95.2	50.1	70.3	55.8	67.8	
w4a16 rollout GSPO	22.0/18.7	55.4	84.0	47.4	47.1	45.7	89.3	42.4	65.3	47.9	61.2	
w4a16 QaRL TBPO	27.5/22.0	62.9	87.2	51.4	56.1	51.2	96.6	48.2	68.0	55.4	67.05	

Table 1: Main results on in-distribution math and out-of-distribution benchmarks. LiveCodeBench results are reported in *pass@4*, while all other metrics use *pass@1*.

than the BF16 RL baseline, whereas QaRL TBPO exhibits markedly improved stability, converging to BF16 comparable rewards. For instance, on the Qwen3-8B model, the average in-distribution math performance drops significantly from 51.8% (BF16) to 43.9% with quantized rollout training, while QaRL TBPO successfully maintain the performance to 48.9%, only 2.9% drop. Critically, QaRL TBPO recovers most of the degradation from quantized-rollout training, achieving near-baseline performance. Although quantized rollout training under GSPO still suffers from router shift at each forward pass, TBPO remains stable on MoE models, achieving an average math score of 51.2%, nearly matching the 52.1% baseline.

Beyond math, QaRL TBPO improves OOD performance vs quantized rollout training while matching BF16, demonstrating that gains arise

from stable optimization and better generalization rather than overfitting.

## 4.2 Ablation

**Optimize objectives** We conduct ablations to validate TBPO’s effectiveness in mitigating error-token interference under QaRL (Fig. 8). **Overall:** TBPO achieves stable learning, high reward, and tight KL control, while alternatives suffer from KL drift or low sample efficiency. **GSPO** is unstable error tokens in negatives cause KL drift and collapse, degrading reward. **MIS GSPO:** Rejection sampling reduces data efficiency and limits reward ceiling. **Positive GRPO:** Discarding negatives loses exploration and limits performance. **On-policy GRPO:** Single updates per batch reduce error amplification but hurt efficiency. **Dual-clip GRPO:** Clipping bounds extreme ratios, yet to-

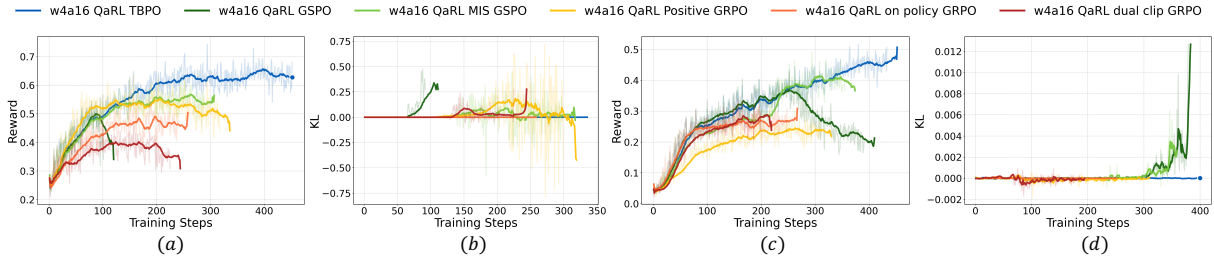


Figure 8: Training dynamics (Reward/KL) of Qwen2.5-Math 1.5B (a-b) and 7B (c-d) across different optimization objectives.

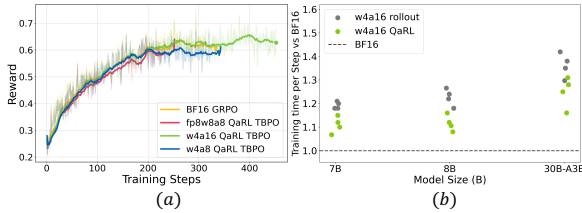


Figure 9: (a) Reward curves under different quant scheme. (b) Per-step training time speedup ratio.

kens after error tokens remain contaminated, causing incorrect learning. For comprehensive analysis of SAPO, please refer to Appendix D.4.

**Quantization Scheme** We further ablate quantization schemes in QaRL-TBPO, comparing FP8W8A8, W4A16, and W4A8 against the BF16 GRPO baseline (Fig. 9(a)). No substantial differences emerge: all bit-widths exhibit similar reward curves and nearly identical final rewards. This suggests that, once stabilized via TBPO, RL convergence is largely insensitive to the specific low-bit format, and performance gains are not tied to a particular quantization choice. We adopt W4A16 for most experiments due to its superior speed on larger models and broad hardware compatibility, whereas W4A8 requires more complex kernel support. Results of TBPO used on quantized rollout training are detailed in appendix D.2.

**Speed** Fig. 9(b) reports per-step training time normalized to BF16. Across 7B, 8B, and 30B-A3B MoE models, both QaRL and quantized rollout training achieve speedups over BF16. For large-scale MoE models, quantized rollout training delivers a  $1.4\times$  speedup, while QaRL achieves a  $1.3\times$  speedup, since QaRL incurs modest overhead from in-training quantization. Comparison of different quantization schemes are provided in appendix 11. We opted not to implement FP8 KV quantization, as KV quant currently fails to provide meaningful throughput benefits in vLLM.

## 5 Related Work

**Quantization.** Quantization accelerates and compresses LLMs via post-training quantization (PTQ) or quantization-aware training (QAT). PTQ methods such as GPTQ (Frantar et al., 2022) and AWQ (Lin et al., 2024) avoid retraining but often struggle on reasoning models, where long CoT generation induces distribution shift and error accumulation beyond static calibration. QAT methods (e.g., LLM-QAT (Liu et al., 2024), EfficientQAT (Chen et al., 2025)) train with quantization noise for better robustness, while fully quantized training (Wang et al., 2025b) executes low-bit arithmetic end-to-end to further improve robustness and throughput.

**Reinforcement Learning in LLMs.** Modern reasoning models (e.g., DeepSeek-R1, Qwen3 (Yang et al., 2025), Kimi-K2 (Moonshot AI, 2025)) commonly build on GRPO. Recent variants include DAPO (higher-bound clipping to promote exploration), GSPO (sequence-level importance sampling for MoE stability), and ASPO (Wang et al., 2025a) (asymmetric sampling for low-probability tokens). Building on this line, we focus on stabilizing RL under quantized rollouts and QAT by suppressing quantization-induced errors during generation.

## 6 Conclusion and Limitations

We propose QaRL to mitigate the severe training-inference mismatch in quantized-rollout RL via rollout-aligned training, and identify error tokens as a key driver of collapse under quantized rollouts, addressed by TBPO with a sequence-level objective and dual clipping on negative samples. Looking forward, we plan to explore fully quantized RL training, and to replace costly, low-utilization sequence-level optimization with token-level approximations that retain stability while improving efficiency and sample usage.

## References

- 535 Mengzhao Chen, Wenqi Shao, Peng Xu, Jiahao Wang,  
536 Peng Gao, Kaipeng Zhang, and Ping Luo. 2025. Ef-  
537 ficientqat: Efficient quantization-aware training for  
538 large language models. In *Proceedings of the 63rd  
539 Annual Meeting of the Association for Computa-  
540 tional Linguistics (Volume 1: Long Papers)*, pages  
541 10081–10100.
- 542 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,  
543 Ashish Sabharwal, Carissa Schoenick, and Oyvind  
544 Tafjord. 2018. Think you have solved question an-  
545 swering? try arc, the ai2 reasoning challenge. *arXiv  
546 preprint arXiv:1803.05457*.
- 547 Elias Frantar, Saleh Ashkboos, Torsten Hoeffler, and  
548 Dan Alistarh. 2022. Gptq: Accurate post-training  
549 quantization for generative pre-trained transformers.  
550 *arXiv preprint arXiv:2210.17323*.
- 551 Chang Gao, Chujie Zheng, Xiong-Hui Chen, Kai  
552 Dang, Shixuan Liu, Bowen Yu, An Yang, Shuai  
553 Bai, Jingren Zhou, and Junyang Lin. 2025. Soft  
554 adaptive policy optimization. *arXiv preprint  
555 arXiv:2511.20347*.
- 556 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,  
557 Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong  
558 Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.  
559 Deepseek-r1: Incentivizing reasoning capability in  
560 llms via reinforcement learning. *arXiv preprint  
561 arXiv:2501.12948*.
- 562 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding  
563 Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han,  
564 Yujie Huang, Yuxiang Zhang, and 1 others. 2024.  
565 Olympiadbench: A challenging benchmark for pro-  
566 moting agi with olympiad-level bilingual multi-  
567 modal scientific problems. In *Proceedings of the  
568 62nd Annual Meeting of the Association for Compu-  
569 tational Linguistics (Volume 1: Long Papers)*, pages  
570 3828–3850.
- 571 Horace He and Thinking Machines Lab. 2025.  
572 [Defeating nondeterminism in llm inference](#).  
573 *Thinking Machines Lab: Connectionism*.  
574 [https://thinkingmachines.ai/blog/defeating-  
575 nondeterminism-in-llm-inference/](https://thinkingmachines.ai/blog/defeating-nondeterminism-in-llm-inference/).
- 576 Jacob Hilton, Karl Cobbe, and John Schulman. 2022.  
577 Batch size-invariance for policy optimization. *Ad-  
578 vances in Neural Information Processing Systems*,  
579 35:17086–17098.
- 580 Jian Hu, Xibin Wu, Wei Shen, Jason Klein Liu, Weixun  
581 Wang, Songlin Jiang, Haoran Wang, Hao Chen, Bin  
582 Chen, Wenkai Fang, Xianyu, Yu Cao, Haotian Xu,  
583 and Yiming Liu. 2025. [OpenRLHF: A ray-based  
584 easy-to-use, scalable and high-performance RLHF  
585 framework](#). In *Proceedings of the 2025 Conference  
586 on Empirical Methods in Natural Language Pro-  
587 cessing: System Demonstrations*, pages 656–666,  
588 Suzhou, China. Association for Computational Lin-  
589 guistics.
- Wei Huang, Yi Ge, Shuai Yang, Yicheng Xiao, 590  
Huizi Mao, Yujun Lin, Hanrong Ye, Sifei Liu, 591  
Ka Chun Cheung, Hongxu Yin, and 1 others. 2025. 592  
[Qerl: Beyond efficiency–quantization-enhanced re-  
593 inforcement learning for llms](#). *arXiv preprint  
594 arXiv:2510.11696*. 595
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richard- 596  
son, Ahmed El-Kishky, Aiden Low, Alec Helyar, 597  
Aleksander Madry, Alex Beutel, Alex Carney, and 1 598  
others. 2024. Openai o1 system card. *arXiv preprint  
599 arXiv:2412.16720*. 600
- Keller Jordan, Yuchen Jin, Vlado Boza, You Jiacheng, 601  
Franz Cesista, Laker Newhouse, and Jeremy Bern- 602  
stein. 2024. [Muon: An optimizer for hidden layers  
603 in neural networks](#). 604
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying 605  
Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. 606  
Gonzalez, Hao Zhang, and Ion Stoica. 2023. Effi- 607  
cient memory management for large language model 608  
serving with pagedattention. In *Proceedings of the  
609 ACM SIGOPS 29th Symposium on Operating Sys-  
610 tems Principles*. 611
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri 612  
Edwards, Bowen Baker, Teddy Lee, Jan Leike, 613  
John Schulman, Ilya Sutskever, and Karl Cobbe. 614  
2023. Let’s verify step by step. *arXiv preprint  
615 arXiv:2305.20050*. 616
- Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, 617  
Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, 618  
Xingyu Dang, Chuang Gan, and Song Han. 2024. 619  
[Awq: Activation-aware weight quantization for on-  
620 device llm compression and acceleration](#). *Proceed-  
621 ings of machine learning and systems*, 6:87–100. 622
- Jiacai Liu, Yingru Li, Yuqian Fu, Jiawei Wang, Qian 623  
Liu, and Yu Shen. 2025a. [When speed kills sta-  
624 bility: Demystifying RL collapse from the training-  
625 inference mismatch](#). 626
- Liyuan Liu, Feng Yao, Dinghui Zhang, Chengyu 627  
Dong, Jingbo Shang, and Jianfeng Gao. 2025b. 628  
[Flashrl: 8bit rollouts, full power rl](#). 629
- Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie 630  
Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, 631  
Raghuraman Krishnamoorthi, and Vikas Chandra. 632  
2024. Llm-qat: Data-free quantization aware train- 633  
ing for large language models. In *Findings of the As-  
634 sociation for Computational Linguistics: ACL 2024*,  
635 pages 467–484. 636
- Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, 637  
William Y. Tang, Manan Roongta, Colin Cai, Jeffrey 638  
Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica. 639  
2025. Deepscaler: Surpassing o1-preview with a 640  
1.5b model by scaling rl. [https://pretty-radio-  
641 -b75.notion.site/DeepScaler-Surpassing-O  
642 1-Preview-with-a-1-5B-Model-by-Scaling-R  
643 L-19681902c1468005bed8ca303013a4e2](https://pretty-radio-b75.notion.site/DeepScaler-Surpassing-O1-Preview-with-a-1-5B-Model-by-Scaling-RL-19681902c1468005bed8ca303013a4e2). Notion  
644 Blog. 645

646	Moonshot AI. 2025. Kimi-k2: Thinking and reasoning. <a href="https://moonshotai.github.io/Kimi-K2/thinking.html">https://moonshotai.github.io/Kimi-K2/thinking.html</a> . Accessed: 2025-12-22.	
647		
648		
649	Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori B Hashimoto. 2025. s1: Simple test-time scaling. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 20286–20332.	
650		
651		
652		
653		
654		
655		
656	Alex Gu Wen-Ding Li Fanjia Yan Tianjun Zhang Sida Wang Armando Solar-Lezama Koushik Sen Ion Stoica Naman Jain, King Han. 2024. Livecodebench: Holistic and contamination free evaluation of large language models for code. <i>arXiv preprint</i> .	
657		
658		
659		
660		
661	Penghui Qi, Zichen Liu, Xiangxin Zhou, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. 2025. Defeating the training-inference mismatch via fp16. <i>arXiv preprint arXiv:2510.26788</i> .	
662		
663		
664		
665	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 <i>First Conference on Language Modeling</i> .	
666		
667		
668		
669		
670	John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. 2015. Trust region policy optimization. In <i>International conference on machine learning</i> , pages 1889–1897. PMLR.	
671		
672		
673		
674	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> .	
675		
676		
677		
678	SGLang Team. 2025. <i>Deterministic Inference</i> . Accessed: 2025-12-21.	
679		
680	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, and 1 others. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. <i>arXiv preprint arXiv:2402.03300</i> .	
681		
682		
683		
684		
685		
686	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2024. Hybridflow: A flexible and efficient rlhf framework. <i>arXiv preprint arXiv:2409.19256</i> .	
687		
688		
689		
690		
691	Jiakang Wang, Runze Liu, Lei Lin, Wenping Hu, Xiu Li, Fuzheng Zhang, Guorui Zhou, and Kun Gai. 2025a. Asp0: Asymmetric importance sampling policy optimization. <i>arXiv preprint arXiv:2510.06062</i> .	
692		
693		
694		
695		
696	Wenjun Wang, Shuo Cai, Congkai Xie, Mingfa Feng, Yiming Zhang, Zhen Li, Kejing Yang, Ming Li, Jiannong Cao, and Hongxia Yang. 2025b. In-fir2: A comprehensive fp8 training recipe for reasoning-enhanced language models. <i>arXiv preprint arXiv:2509.22536</i> .	
697		
698		
699		
700		
701		
	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, and 1 others. 2024. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. <i>Advances in Neural Information Processing Systems</i> , 37:95266–95290.	702 703 704 705 706 707 708
	Bram Wasti, Wentao Ye, Teja Rao, Michael Goin, Paul Zhang, Tianyu Liu, Natalia Gimelshein, Woosuk Kwon, Kaichao You, and Zhuohan Li. 2025. No more train-inference mismatch: Bitwise consistent on-policy reinforcement learning with vllm and torchtitan. Accessed: 2025-12-21.	709 710 711 712 713 714
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	715 716 717 718 719 720
	Haocheng Xi, Han Cai, Ligeng Zhu, Yao Lu, Kurt Keutzer, Jianfei Chen, and Song Han. 2024. Coat: Compressing optimizer states and activation for memory-efficient fp8 training. <i>arXiv preprint arXiv:2410.19313</i> .	721 722 723 724 725
	Jianhao Yan, Yafu Li, Zican Hu, Zhi Wang, Ganqu Cui, Xiaoye Qu, Yu Cheng, and Yue Zhang. 2025. Learning to reason under off-policy guidance. <i>Preprint</i> , arXiv:2504.14945.	726 727 728 729
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	730 731 732 733 734
	Feng Yao, Liyuan Liu, Dinghui Zhang, Chengyu Dong, Jingbo Shang, and Jianfeng Gao. 2025. Your efficient rl framework secretly brings you off-policy rl training.	735 736 737 738
	Deheng Ye, Zhao Liu, Mingfei Sun, Bei Shi, Peilin Zhao, Hao Wu, Hongsheng Yu, Shaojie Yang, Xipeng Wu, Qingwei Guo, and 1 others. 2020. Mastering complex control in moba games with deep reinforcement learning. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 34, pages 6672–6679.	739 740 741 742 743 744 745
	Yixin Ye, Yang Xiao, Tiantian Mi, and Pengfei Liu. 2025. Aime-preview: A rigorous and immediate evaluation framework for advanced mathematical reasoning. <a href="https://github.com/GAIR-NLP/AIME-Preview">https://github.com/GAIR-NLP/AIME-Preview</a> . GitHub repository.	746 747 748 749 750
	Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, and 1 others. 2025. Dapo: An open-source llm reinforcement learning system at scale. <i>arXiv preprint arXiv:2503.14476</i> .	751 752 753 754 755

756	Ziyang Zhang, Xinheng Ding, Jiayi Yuan, Rixin Liu,	<b>Appdenix</b>	774
757	Huizi Mao, Jiarong Xing, and Zirui Liu. 2025. <a href="#">De-</a>		
758	<a href="#">terministic inference across tensor parallel sizes that</a>	<b>A The Use of Large Language Models</b>	775
759	<a href="#">eliminates training-inference mismatch.</a> <i>Preprint,</i>	<b>(LLMs)</b>	776
760	<a href="#">arXiv:2511.17826.</a>		
761	Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai	A large language model was utilized for grammat-	777
762	Dai, Huazuo Gao, Jiashi Li, Liyue Zhang, Panpan	ical and stylistic refinement of the manuscript. Its	778
763	Huang, Shangyan Zhou, Shirong Ma, and 1 oth-	role was strictly limited to text editing and polish-	779
764	ers. 2025. Insights into deepseek-v3: Scaling chal-	ing to enhance clarity. All research ideas, experi-	780
765	lenges and reflections on hardware for ai architec-	mental design, and analytical content are the origi-	781
766	tures. In <i>Proceedings of the 52nd Annual Interna-</i>	nal work of the authors.	782
767	<i>tional Symposium on Computer Architecture</i> , pages		
768	1731–1745.		
769	Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui	<b>B Experiment Setting</b>	783
770	Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong		
771	Liu, Rui Men, An Yang, and 1 others. 2025.	<b>Training.</b> We use Ver1 (Sheng et al., 2024) as	784
772	Group sequence policy optimization. <i>arXiv preprint</i>	the training framework and vLLM (Kwon et al.,	785
773	<i>arXiv:2507.18071.</i>	2023) as the inference engine. All experiments are	786
		conducted on $8 \times$ NVIDIA H800 GPUs.	787
		<b>Evaluation Dataset.</b> We evaluate our approach	788
		on standard in distribution mathematical datasets	789
		including AIME2024, AIME2025 (Ye et al.,	790
		2025), AMC, Math-500 (Lightman et al., 2023),	791
		OlympiadBench (He et al., 2024) and Min-	792
		erva. To further investigate the generalizabil-	793
		ity of quantized training, we extend evaluation	794
		to out-of-distribution benchmarks such a out-of-	795
		distribution benchmarks ARC-Challenge (Clark	796
		et al., 2018), GPQA-Diamond (Rein et al., 2024),	797
		LiveCodeBench (Naman Jain, 2024), MMLU	798
		Pro (Wang et al., 2024).	799
		<b>Hyperparameters.</b>	800
		• seq_clip_ratio_high is $\epsilon_h$ ,	801
		• neg_seq_clip_ratio_high is $\delta_h$ and	802
		neg_seq_clip_ratio_low is $\delta_l$ ,	803
		• seq_tis_imp_ratio_cap is $c$ in equation 5.	804
		<b>C Batch invariant kernel</b>	805
		The root cause of output inconsistency across dif-	806
		ferent kernel implementations lies in the finite pre-	807
		cision of floating point accumulators. Due to the	808
		non associativity of floating point addition, where	809
		$(a+b)+c = a+(b+c)$ the final result is sensitive to	810
		the order of operations. To ensure bit-wise repro-	811
		ducibility, it is necessary to enforce a fixed sum-	812
		mation order regardless of the batch configuration,	813
		leading to the design of a Batch-Invariant Kernel.	814
		However, imposing a strict execution order often	815
		comes at the cost of performance, as it restricts the	816
		asynchronous parallelism and dynamic scheduling	817
		inherent to GPU architectures.	818

Parameter Name	Value
trainer.nnodes	1
trainer.n_gpu_per_node	8
data.train_batch_size	512
data.max_prompt_length	2048
data.max_response_length	16384
rollout.n	8
rollout.temperature	1.0
rollout.top_p	1.0
val_kwargs.temperature	0.6
actor.ppo_mini_batch_size	64
actor.ppo_max_token_len_per_gpu	22528
optim.opt_type	Muon
optim.lr	1e-6
optim.weight_decay	0.01
actor.use_kl_loss	False
actor.seq_clip_ratio_high	0.0004
actor.seq_clip_ratio_low	0.0003
actor.neg_seq_clip_ratio_high	0.0007
actor.neg_seq_clip_ratio_low	0.0003
actor.seq_tis_imp_ratio_cap	2

Table 2: Hyperparameters for Experiment.

## D More Experimental Results

### D.1 Speed on different quant scheme

Fig. 11 illustrates the per-step training latency of various quantization schemes on Qwen3-30B-A3B (MoE), normalized to the BF16 baseline (dashed line at 1.0). Our results show that efficiency gains become increasingly significant from W8 to W4. This trend underscores that MoE training is primarily memory/IO-bound; since MoE operators are almost inherently memory-bound during decoding, the weight bit-width directly dictates computational efficiency. Furthermore, lower precision reduces the memory footprint per model instance, minimizing the required GPU count and alleviating inter-GPU communication overheads that becomes even more pronounced when intra-node interconnects (e.g., NVLink) are fully utilized. Consequently, we adopt W4 as the default configuration for our primary experiments.

### D.2 More Ablation on TBPO

To further evaluate the effectiveness of TBPO, particularly its capability to mitigate the negative impact of "error tokens" generated by quantized rollout engines, we conduct an ablation study against GRPO under different precision settings. As

Method	Qwen3-8B-Base			
	MATH-500	AIME25	AMC	Avg.
BF16 GRPO	88.1	19.3	64.3	57.2
BF16 TBPO	86.3	20.0	62.4	56.2
w4a16 GRPO	82.2	12.5	53.0	49.2
w4a16 TBPO	84.9	13.3	58.3	52.1

Table 3: Performance comparison between GRPO and TBPO under full-precision (BF16) and quantized (w4a16) rollout training settings.

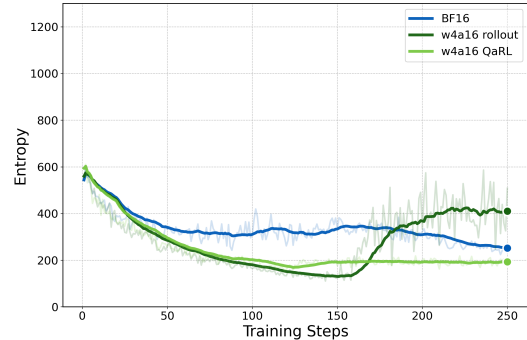


Figure 10: Comparison of RL training entropy

demonstrated in Table 3, TBPO achieves performance on par with GRPO under BF16 precision. However, under w4a16 quantized rollout training, TBPO exhibits markedly superior robustness by effectively neutralizing error token interference during optimization, thereby delivering an average margin of 2.9 points over GRPO. While TBPO provides some degree of mitigation for the mismatch issues arising from quantized rollout training, the fundamental resolution lies in leveraging QAT/QaRL or fully quantized training regimes to eliminate systemic inconsistencies at their core.

### D.3 Entropy of quantized RL training

Whereas Huang et al. (2025) Fig. 5 posits that quantization errors during training may attenuate

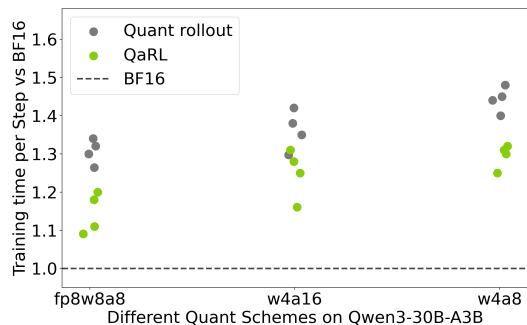


Figure 11: Different quant scheme speed on Qwen3-30B-A3B MoE model.

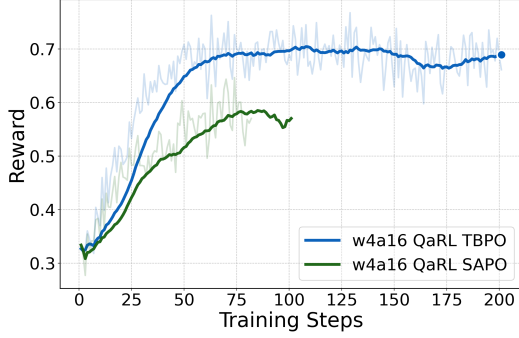


Figure 12: Comparison of SAPO on QaRL

the entropy reduction in RL, our empirical findings present a more nuanced picture. Across both quantized rollout training and QAT/QaRL paradigms (Fig. 10), quantization discrepancies are progressively assimilated throughout the optimization process, yielding no discernible entropy elevation relative to the BF16 baseline. The entropy escalation observed in the latter stages of w4a16 quantized rollout training stems from the over-optimization of error tokens, which precipitates repetitive generation patterns.

#### D.4 Comparison of soft weighted methods SAPO

SAPO (Gao et al., 2025) introduces a soft adaptive weighting mechanism as a nuanced alternative to hard clipping, aiming to bolster optimization stability in MoE architectures. Through dynamic reweighting of token contributions assigning diminishing importance to tokens whose advantage ratios deviate further from unity SAPO adeptly discriminates between erroneous tokens (demanding substantial correction) and exploration tokens (that foster policy discovery). Nevertheless, as evidenced in Fig. 12, SAPO remains unable to surpass TBPO in the context of quantized RL training. We posit that this suboptimal performance stems from SAPO’s continued assimilation of responses containing error tokens, despite attenuating their individual contributions to policy updates. Critically, such response sequences are globally off-distribution, rendering them intrinsically unsuitable for stable policy learning. This observation reaffirms the paramount importance of constraining sequence-level optimization objectives within the trust region, while simultaneously attesting to the elegance and principled simplicity inherent in TBPO’s design.

## E Detail of TBPO

To effectively optimize the policy while maintaining training stability, we define the sequence-level importance weights and mismatch factors. Specifically, for a given query  $q$  and its corresponding output sequence  $o$  of length  $L$ , we introduce the sequence-level proximity ratio  $r_{\text{seq-prox}}(\theta)$  and the sequence-level mismatch weight  $w_{\text{seq-mismatch}}(\theta)$  as illustrated in Fig. 13.

Instead of using a simple product of token-level probabilities, which often leads to vanishing or exploding gradients in long-context scenarios, we employ the geometric mean of token-level ratios (formulated as the exponential of the average log-difference). This design ensures that the importance weights remain numerically stable across varying sequence lengths.  $r_{\text{seq-prox}}(\theta)$  serves to constrain the policy update within a reliable trust region by measuring the drift from the previous learner, while  $w_{\text{seq-mismatch}}(\theta)$  accounts for the distributional shift between the historical sampling policy and the current optimization baseline.

