GEPO: GROUP EXPECTATION POLICY OPTIMIZATION FOR STABLE HETEROGENEOUS REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

Abstract

As single-center computing approaches power constraints, decentralized training becomes essential. However, traditional Reinforcement Learning (RL) methods, crucial for enhancing large model post-training, cannot adapt to decentralized distributed training due to the tight coupling between parameter learning and rollout sampling. For this, we propose HeteroRL, a heterogeneous RL architecture that decouples these processes, enabling stable training across geographically distributed nodes connected via the Internet. The core component is Group Expectation Policy Optimization (GEPO), an asynchronous RL algorithm robust to latency caused by network delays or heterogeneity in computational resources. Our study reveals that high latency significantly increases KL divergence, leading to higher variance of importance weights and training instability. GEPO mitigates this issue by using group expectation weighting to exponentially reduce the variance of importance weights, with theoretical guarantees. Experiments show GEPO achieves superior stability—only a 3% performance drop from online to 1800s latency—and reduces the best-to-last gap by 85% versus GSPO (Δ =1.8 vs. 12.0) while attaining the highest scores, highlighting its effectiveness in decentralized, resource-heterogeneous environments.

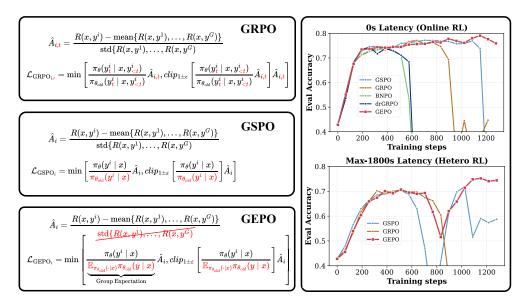


Figure 1: Left: GEPO improves upon GRPO and GSPO by employing group-level importance weights to enhance training stability. Right: In both zero-delay (online) and high-delay (up to 1800 seconds) heterogeneous reinforcement learning scenarios, GEPO demonstrates superior stability and better evaluation performance.

1 Introduction

Training ever-larger AI models (Achiam et al., 2023; Dubey et al., 2024; Yang et al., 2025) is pushing the limits of single datacenters, making decentralized training across geographically distributed, heterogeneous nodes connected via the Internet an increasingly necessary pursuit (Team et al., 2025; Noukhovitch et al., 2024). Reinforcement Learning (RL), crucial for post-training LLMs on complex tasks like mathematical reasoning (Shao et al., 2024), faces a fundamental systemic challenge in this emerging paradigm: traditional RL frameworks (Guo et al., 2025; Stiennon et al., 2020; Bai et al., 2022; Wu et al., 2025; Dai et al., 2024; Fu et al., 2025) are architecturally incompatible with decentralized environments. Their tight coupling between rollout sampling and parameter learning demands strict synchronization—a requirement that becomes untenable under the high network latency and computational heterogeneity inherent in real-world distributed settings.

This architectural incompatibility manifests in two critical bottlenecks. First, synchronous frameworks force computational resources (e.g., GPUs) to idle while waiting for the slowest processes—such as generating long reasoning chains—severely constraining efficiency (Fu et al., 2025). Second, and more fundamentally, the inevitable network latency inherent in decentralized, Internet-connected environments creates a temporal gap (policy staleness) between the sampler (generating data) and the learner (updating parameters). Most existing RL algorithms, designed for homogeneous, low-latency clusters, are ill-equipped to handle this staleness. As our analysis reveals, high latency significantly inflates KL divergence, causing the variance of importance weights to explode—ultimately leading to training instability or reward collapse (Song et al., 2023). This renders conventional RL methods impractical for real-world, geographically distributed training scenarios.

To tackle these systemic bottlenecks, we introduce **HeteroRL** (Heterogeneous Reinforcement Learning), a novel RL framework explicitly architected for asynchronous, geographically distributed, and resource-heterogeneous environments. HeteroRL is designed to enable efficient and stable training of large language models for complex tasks such as mathematical reasoning, even under high network latency. At its core, HeteroRL decouples the two computationally intensive phases of the RL pipeline — rollout sampling and parameter learning — and deploys them on physically or logically independent nodes with potentially heterogeneous hardware (e.g., mixing NVIDIA and Ascend chips). The sampler nodes continuously generate reasoning trajectories without interruption, while the learner node asynchronously consumes this data to update model parameters. Critically, neither component waits for the other: communication occurs infrequently and tolerates high latency, with model checkpoints and rollout batches exchanged over the Internet.

To address the instability arising from KL divergence under high-latency conditions, we introduce **Group Expectation Policy Optimization (GEPO)**, a novel policy gradient algorithm that stabilizes asynchronous RL under high latency by replacing fragile token/sample-level importance weights with robust group-level importance weights — allowing samplers and learners to operate independently, communicating infrequently and tolerating arbitrary delays. This shift fundamentally improves the quality of gradient estimation — transforming a high-variance, unstable estimator into a low-variance, robust one, especially under large policy divergence. As we prove in Theorem 1, GEPO exponentially reduces the variance of importance weights under high KL divergence — precisely the regime where traditional methods like GRPO and GSPO collapse. Crucially, GEPO is not an ad hoc fix — it is a principled algorithmic response to the root cause of instability: variance explosion under policy divergence.

In summary, our key contributions are as follows:

Framework: We propose HeteroRL, an asynchronous reinforcement learning framework designed for heterogeneous compute networks, enabling decentralized training of large language models (LLMs) on mathematical reasoning tasks.

Insight: We identify a strong correlation between latency and the KL divergence between the rollout sampler and the learner. High latency induces high KL divergence, leading to training instability and reward collapse.

Algorithm: We introduce Group Expectation Policy Optimization (GEPO), which improves upon the importance sampling mechanism in GRPO (Shao et al., 2024). We theoretically show that GEPO exponentially reduces the variance of importance weights, and empirically demonstrate its superior stability and efficiency — not only under high-latency conditions, but also in the ideal zero-latency setting.

This work provides both algorithmic and system-level advancements for scalable LLM RL-training and establishes a practical foundation for large-scale distributed AI training in future heterogeneous compute network environments.

2 Background

2.1 Problem Definition and Notation

Consider a standard policy gradient framework. Let π_{θ} denote the language model policy (i.e., the Actor) parameterized by θ , x be an input prompt of a Dataset \mathcal{D} (e.g., math problems), and y be the output sequence generated by the model (e.g., a chain-of-thought solution). We define the following core notation:

- π_{θ_k} (short for q): the policy used by the *sampler* at time step k to generate rollout trajectories.
- $\pi_{\theta_{k+\tau}}$ (short for p): the latest policy at the learner at time step $k+\tau$, used for gradient updates.
- $\tau(\geq 0)$: policy staleness, representing the discrepancy in policy versions between the sampler and the learner, caused by network delays and computational asynchrony.
- y: a trajectory sampled from the stale policy π_{θ_k} . y_t^i denotes the t-th token of the i-th response in a group.
- r(x,y): the reward for response y given input x.
- A(x,y): the advantage for response y to input x, typically defined as A(x,y) = r(x,y) b(x), where b(x) is a baseline reward computed for input x. In this paper, we use the within-group average reward (Shao et al., 2024) as the baseline $b(x) = \frac{1}{G} \sum_{i=1}^{G} r(x,y^i)$.

The goal of HeteroRL is to optimize the policy π_{θ} to maximize the expected cumulative reward. To reduce gradient variance, an advantage function is used, leading to the objective:

$$\mathcal{L}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta_k}(\cdot|x)} \left[\frac{\pi_{\theta_k + \boxed{T}}(y|x)}{\pi_{\theta_k}(y|x)} \cdot A(x, y) \right], \tag{1}$$

where τ is a random variable: $\tau \sim p(\text{Sync_Model_Time}, \text{Sync_Data_Time}, \text{Rollout_Time})$. For online RL, $\tau \equiv 0$.

3 GEPO: Group Expectation Policy Optimization

Our method builds upon the group-based policy optimization paradigm of GRPO and introduces the group expectation importance sampling mechanism. We emphasize a paradigm shift from *token-level* to *group-level* importance weighting, which significantly reduces the variance of importance weights and alleviates gradient instability during training.

3.1 Group Expectation Importance Weighting

To enhance the stability of importance weights, we propose the Group Expectation Importance Weight (GEIW), which replaces the individual proposal probability q(y|x) in the standard importance weight $\frac{p(y|x)}{q(y|x)}$ with its group-wise expected value under the current prompt x, denoted as $\widehat{\mathbb{E}}_q[q(y|x)]$. Inspired by GRPO, for each input x, we generate a group of G responses $\{y^1,\ldots,y^G\} \sim q(\cdot|x)$ to form a sampling group. Since G is typically much smaller than the full policy space and top-P/top-K sampling leads to $\sum_{i=1}^G q(y^i|x) \gg 1$,

the vector $(q(y^1|x), \ldots, q(y^G|x))$ does not constitute a valid probability distribution. Simply using the arithmetic mean $\frac{1}{G}\sum_{i=1}^G q(y^i|x)$ would introduce bias due to ignoring the relative sampling probabilities. To obtain a more accurate estimate, we employ a weighted expectation:

$$\widehat{\mathbb{E}}_{q}[q(y|x)] \approx \sum_{i=1}^{G} \widehat{q(y^{i}|x)} \cdot q(y^{i}|x) = \frac{\sum_{i=1}^{G} q(y^{i}|x)^{2}}{\sum_{i=1}^{G} q(y^{i}|x)},$$
(2)

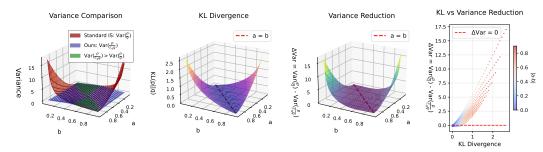
where $\widehat{q(y^i|x)} = \frac{q(y^i|x)}{\sum_{i=1}^G q(y^i|x)}$ is the within-group normalized probability, serving as an empirical estimate of the sampling likelihood of each y_i . We define the GEIW importance weight as:

$$w_{\text{GEIW}}(y|x) = \frac{p(y|x)}{\widehat{\mathbb{E}}_q[q(y|x)]}.$$
 (3)

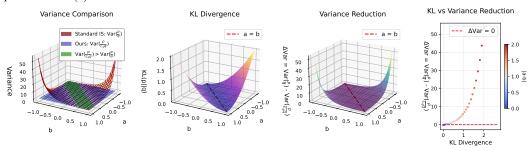
The key advantages of this mechanism are as follows:

Numerically stable and gradient-effective: The denominator is decoupled from any single q(y|x), avoiding extreme weight values when individual proposal probabilities approach zero. Although $\operatorname{clip}(1\pm\epsilon)$ can also improve numerical stability, the gradients of the clipped tensors will be set to zero, effectively skipping this data point (ineffective gradient).

Biased yet low-variance: By leveraging within-group statistical information, GEIW provides a more reliable scale estimate. Even under large divergence between p and q, $\widehat{\mathbb{E}}_q[q(y|x)]$ remains well-conditioned, effectively preventing gradient explosion. Although this estimator introduces a small bias (w_{GEIW} is a biased estimator), both theoretical analysis (see Theorem 1) and empirical results demonstrate that it significantly reduces variance under high KL divergence, yielding more stable gradient directions and improved training convergence.



(a) Variance comparison of $\frac{p}{q}$ and $\frac{p}{\mathbb{E}_q[q]}$ under Bernoulli distributions, where $p \sim \text{Bernoulli}(a)$ and $q \sim \text{Bernoulli}(b)$.



(b) Variance comparison of $\frac{p}{q}$ and $\frac{p}{\mathbb{E}_q[q]}$ under Gaussian distributions, where $p \sim \mathcal{N}(a,1)$ and $q \sim \mathcal{N}(b,1)$.

Figure 2: In high-KL regions, $\operatorname{Var}\left[\frac{p(y|x)}{\widehat{\mathbb{E}}_q[q(y|x)]}\right] \ll \operatorname{Var}\left[\frac{p(y|x)}{q(y|x)}\right]$.

Theorem 1. Let p,q be discrete probability distributions. Then there exists a constant C such that:

 $\operatorname{Var}\left[\frac{p(y|x)}{q(y|x)}\right] - \operatorname{Var}\left[\frac{p(y|x)}{\widehat{\mathbb{E}}_{q}[q(y|x)]}\right] \ge \left[\exp\left(D_{\mathrm{KL}}(p\|q)\right)\right] - C. \tag{4}$

In particular, when $D_{\mathrm{KL}}(p||q) > \log C$, it holds that $\mathrm{Var}\left[\frac{p(y|x)}{q(y|x)}\right] > \mathrm{Var}\left[\frac{p(y|x)}{\widehat{\mathbb{E}}_{q}[q(y|x)]}\right]$.

Theorem 1 shows that GEPO can exponentially reduce the variance of importance weights, making it particularly well-suited for heterogeneous RL training under high KL divergence. The full mathematical proof is provided in Appendix A. As shown in Figure 2, we visualize the relationship between KL divergence and importance weight variance when both p and q are Bernoulli or Gaussian distributions with varying parameters. The results indicate that in the high-KL regime, the group expectation approach significantly reduces the variance of importance weights, which benefits training stability under high network latency. Nevertheless, there exist regimes—such as the green regions in the plots—where our method incurs a slight increase in variance.

The difference across all GRPO-like algorithms lies in the computation of the importance weights, as detailed in Listing 1:

```
if self.loss_type in ["grpo","dr_grpo","bnpo"]: # Token level
    coef_1 = learner_token_p / sampler_token_p

elif self.loss_type == "gspo": # Sequence level
    coef_1 = learner_seq_p / sampler_seq_p

elif self.loss_type == "gepo": # Group level
    normalized_q = sampler_seq_p.detach() / (sampler_seq_p.sum().
    detach())

coef_1 = learner_seq_p / (normalized_q * sampler_seq_p).sum()
```

Listing 1: Coefficient computation for different policy optimization methods

3.2 Gradient Comparison Across Tokens

What does the GEPO update do? For a mechanistic understanding of GEPO, it is useful to analyze the gradient of the loss function \mathcal{L}_{GEPO} . The equivalent gradient of each token in a group with respect to the parameters θ of GRPO, GSPO and GEPO can be written as:

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbf{A} \odot \underbrace{\begin{bmatrix} \frac{p'_{1,1}(\boldsymbol{\theta})}{q_{1,1}} & \dots & \frac{p'_{1,T}(\boldsymbol{\theta})}{q_{1,T}} \\ \vdots & \ddots & \vdots \\ \frac{p'_{G,1}(\boldsymbol{\theta})}{q_{G,1}} & \dots & \frac{p'_{G,T}(\boldsymbol{\theta})}{q_{G,T}} \end{bmatrix}}_{\mathbf{GRPO}} \text{ or } \underbrace{\begin{bmatrix} \frac{p'_{1,1}(\boldsymbol{\theta})}{q_{1}} & \dots & \frac{p'_{1,T}(\boldsymbol{\theta})}{q_{1}} \\ \vdots & \ddots & \vdots \\ \frac{p'_{G,1}(\boldsymbol{\theta})}{q_{G}} & \dots & \frac{p'_{G,T}(\boldsymbol{\theta})}{q_{G}} \end{bmatrix}}_{\mathbf{GSPO}} \text{ or } \underbrace{\begin{bmatrix} \frac{p'_{1,1}(\boldsymbol{\theta})}{\mathbb{E}_{q}q} & \dots & \frac{p'_{1,T}(\boldsymbol{\theta})}{\mathbb{E}_{q}q} \\ \vdots & \ddots & \vdots \\ \frac{p'_{G,1}(\boldsymbol{\theta})}{\mathbb{E}_{q}q} & \dots & \frac{p'_{G,T}(\boldsymbol{\theta})}{\mathbb{E}_{q}q} \end{bmatrix}}_{\mathbf{GEPO} \text{ (ours)}},$$

where $\mathbf{A} \in \mathbb{R}^{G \times T}$ is token-level advantages matrix, \odot denotes Hadamard product, $q_{i,t} = q(y_t^i \mid x^i, y_{< t}^i)$, $q_i = q(y^i \mid x)$, and $\mathbb{E}_q q = \widehat{\mathbb{E}}_q[q(y|x)]$. From the perspective of gradient stability, GSPO uses a shared denominator $q(y^i \mid x)$ for all tokens in sequence i, while GEPO further aggregates across the entire group by using a common denominator $\mathbb{E}_q q$. This progression—from token-level (GRPO) to sequence-level (GSPO) to group-level (GEPO) coefficients—demonstrates that coarser importance-weight granularity significantly reduces gradient variance. Empirically, leveraging group-level statistics enhances robustness and stabilizes training, especially under high policy divergence.

4 Experiments

4.1 Experimental Setup

Model, Dataset and Benchmarks We conduct reinforcement learning training and evaluation on the Qwen3-1.7B/8B model. The models are trained by strong-to-weak distillation (Yang et al., 2025), but have not been tuned with any RL. We train the model on 8,290 samples from the MATH level 3–5 dataset (Zeng et al., 2025) and evaluate it by

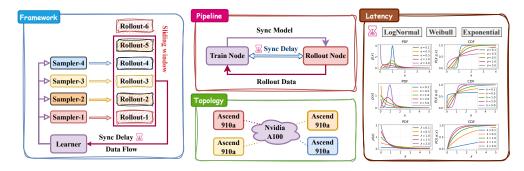


Figure 3: The Overview of HeteroRL. By decoupling sampling and training, HeteroRL enables decentralized distributed RL training of LLMs across five compute nodes: one parameter update node (learner) and four data generation nodes (sampler), forming a star-shaped network topology. Network delays between the sampler and learner nodes are explicitly modeled and can be simulated using stochastic distributions such as the log-normal or Weibull distribution.

reporting average Pass@1 over 8 sampled responses on the MATH500 (Hendrycks et al., 2021), AMC23 (Li et al., 2024), AIME24 (AIME, 2024), and AIME25 (AIME, 2025) benchmarks. To better evaluate the inherent stability of policy optimization algorithms, we remove KL divergence constraints during training under online RL scenario, and use the same KL coefficient under the heterogeneous RL scenario. We compare our method against baseline methods GRPO (Shao et al., 2024), GSPO (Zheng et al., 2025), BNPO (Xiao et al., 2025) and Dr.GRPO (Liu et al., 2025) under both zero-delay and high-delay settings. More experimental details can be found in Section B.1.

Heterogeneous Computing Environment As shown in Figure 3, we perform heterogeneous training across five compute nodes: one learner node and four sampler nodes, forming a star-shaped topology centered at the learner. During training, the sampler nodes generate rollout data, which is transmitted over the network to the learner node in a streaming fashion. The learner updates the model parameters and periodically broadcasts the updated weights back to the sampler nodes. The learner processes incoming rollouts in the order they arrive, operating within a fixed time window for data eligibility. Since data is transmitted in batch units—each containing text, generation probabilities, and rewards—a maximum delay of 1800 seconds is sufficient for typical network conditions. Within this window, the iteration gap (in terms of gradient updates) between the learner and samplers remains within 64 steps.

4.2 Main Experimental Results

Table 1: Performance comparison using Qwen3-1.7B/8B under Online RL (4k limitaion).

| Method | AMC2023 | | AIME2024 | | AIME2025 | | MATH500 | | Average | |
|-----------------------|---------|------|----------|------|----------|------|---------|------|---------|------|
| | 1.7B | 8B | 1.7B | 8B | 1.7B | 8B | 1.7B | 8B | 1.7B | 8B |
| Qwen3-1.7/8B | 44.6 | 70.6 | 10.9 | 32.4 | 14.0 | 26.1 | 72.4 | 87.1 | 35.5 | 54.1 |
| Max Tolerable Delay 0 | | | | | | | | | | |
| BNPO | 59.4 | 78.8 | 27.7 | 44.1 | 23.4 | 29.3 | 83.7 | 91.4 | 48.6 | 60.9 |
| Dr.GRPO | 61.6 | 77.5 | 24.6 | 41.0 | 22.7 | 27.7 | 82.9 | 91.6 | 48.0 | 59.4 |
| GRPO | 60.9 | 81.3 | 30.9 | 42.6 | 24.2 | 31.3 | 83.7 | 92.0 | 49.9 | 61.8 |
| GSPO | 60.3 | 77.8 | 28.5 | 41.8 | 25.0 | 31.3 | 83.9 | 90.9 | 49.4 | 60.5 |
| GEPO (ours) | 62.2 | 85.6 | 31.6 | 44.1 | 25.8 | 37.5 | 84.7 | 92.6 | 51.1 | 65.0 |

In this section, we compare GEPO with baselines under online RL and Hetero RL settings. The experimental results in Tables 1 and 2 demonstrate that GEPO not only achieves

superior performance but also exhibits exceptional stability across both online and Hetero RL settings. Below, we dissect these findings in depth.

Table 2: Performance comparison using Qwen3-8B under Hetero RL (4k limitation).

| Method | AMC2023 | | AIME2024 | | AIME2025 | | MATH500 | | Average | |
|-------------|----------------------------|------|----------|------|----------------------------|------|-----------------|------|---------|------|
| | $\overline{\mathbf{best}}$ | last | best | last | $\overline{\mathrm{best}}$ | last | \mathbf{best} | last | best | last |
| Qwen3-8B | 70.6 | - | 32.4 | - | 26.1 | - | 87.1 | - | 54.1 | - |
| | Max Tolerable Delay 64 | | | | | | | | | |
| GRPO | 71.6 | 71.6 | 38.7 | 35.9 | 27.3 | 27.0 | 88.8 | 88.8 | 56.6 | 55.8 |
| GSPO | 76.2 | 60.0 | 37.9 | 16.4 | 28.9 | 27.3 | 90.7 | 81.9 | 58.4 | 46.4 |
| GEPO (ours) | 83.4 | 82.8 | 42.6 | 37.5 | 33.2 | 32.0 | 91.3 | 90.9 | 62.6 | 60.8 |

In the **online RL setting** (Table 1), GEPO consistently outperforms all baselines across both model sizes and all benchmarks. On Qwen3-8B, it achieves an average score of **65.0**, surpassing GRPO (61.8) and GSPO (60.5) by 3.2 and 4.1 points, respectively. The gain is most notable on AIME2025 (+6.2 points over GRPO/GSPO, 20% relative improvement). Even on the 1.7B model, GEPO sets a new SOTA, exceeding the best baseline by 1.5 points in average, confirming that group-level importance weighting improves gradient quality even without asynchrony.

In the **Hetero RL setting** (Table 2), GEPO's stability advantage becomes decisive. Both GEPO and GSPO improve over GRPO in best performance (+10.6% and +3.2%, respectively). However, GEPO further surpasses GSPO by 7.2% in accuracy while reducing its best-to-last degradation by 85% versus GSPO (Δ =1.8 vs. 12.0), achieving both higher performance and far greater stability. While GSPO's last scores collapse dramatically, GEPO maintains near-peak performance throughout training.

These results validate GEPO's core design: by replacing token or sequence-level importance weights with group-level expectations, it exponentially reduces importance weight variance under high KL divergence (Theorem 1), enabling stable, scalable decentralized RL. GEPO thus sets a new frontier in both performance and stability across ideal and real-world distributed settings. In the HeteroRL setting, the training process recorded in Figure 4 shows

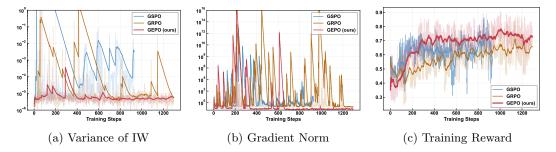


Figure 4: Curves of importance sampling variance, training gradient norm, and train/eval reward under max delay 64. Compared to GRPO and GSPO, GEPO maintains more stable importance weight variance, resulting in less drastic gradient changes, more stable training, and no decline in training reward.

that GRPO stably improves the reward at a slower pace, while GSPO rapidly increases the reward in the first 200 steps but becomes unstable between 500 and 700 steps. As seen in Figure 4a, GEPO exhibits significantly lower variance in importance weights compared to GRPO and GSPO, which experience sharp spikes and fluctuations. These unstable weight variances lead to erratic gradient updates, as evidenced by the large oscillations in gradient norm (Figure 4b) for GRPO and GSPO, especially during early and mid-training phases. In contrast, GEPO's gradient norms remain relatively smooth and bounded, contributing

to stable learning progress. Consequently, the training reward curve (Figure 4c) shows consistent improvement for GEPO without any noticeable decline, whereas GRPO and GSPO exhibit periods of stagnation or even degradation.

4.3 Analysis Experiment

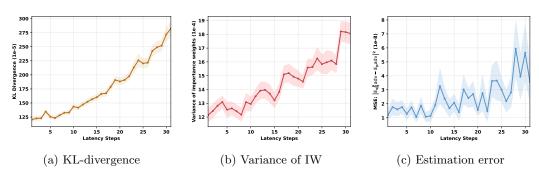
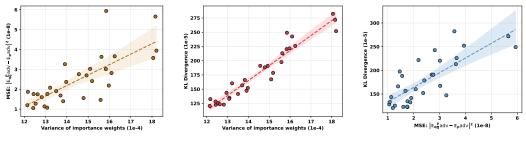


Figure 5: KL-divergence, Variance of IW, and Estimation error are all positively correlated with the number of delay steps.

Impact of Latency As shown in Figure 5, we analyze the changes in KL divergence between the trainer and sampler, variance of importance weights, and estimation error of the expected value of the advantage function (optimization objective) during heterogeneous RL training as latency increases. We observe that latency leads to increased KL divergence (Figure 5a), which in turn causes an increase in the variance of importance weights (Figure 5b), ultimately resulting in increased estimation error of the expected advantage function (Figure 5c). Since the optimization objective is to maximize the estimated expectation of the advantage function, large estimation errors will cause significant fluctuations in gradients, thereby affecting training stability and performance. To show that high latency harms training stability, we compare max delays of 8 and 64 steps. As Figure 7 shows, with 64-step delay—especially near step 900—the KL divergence spikes and evaluation accuracy drops sharply, confirming that latency induces instability. Although GEPO improves stability, it still suffers a performance dip around step 900, highlighting that heterogeneous RL under high latency remains challenging.



(a) The correlation between importance weight variance and estimation error of $\mathbb{E}_p[adv(x,y)]$ is **0.76**.

(b) The correlation between importance weight variance and KL divergence is **0.96**.

(c) The correlation between KL divergence and $\mathbb{E}_p[adv(x,y)]$ estimation error is **0.78**.

Figure 6: Correlation analysis (95% CI) of training delay steps, importance sampling variance, and estimation error of expected advantage function.

Correlation and Causality. Figure 6 quantifies the pairwise correlations among KL-divergence, variance of importance weights, and estimation error of the expected advantage function. The correlation coefficients range from 0.76 to 0.96 ($\alpha=0.05$), confirming a strong statistical association between these variables. This observation empirically supports our hypothesis (illustrated in Figure 5) that increased latency induces higher KL divergence, which in turn amplifies the variance of importance weights and the estimation error, ultimately threatening training stability. However, correlation does not imply strict causation.

While latency is a significant contributing factor to KL divergence, it is not the sole determinant — the model's internal state and the statistical properties of the sampled data also play crucial roles. This explains the observed variance in training outcomes under identical latency: sometimes collapse occurs, sometimes not. Critically, what determines survival versus collapse is not latency itself, but the algorithm's capacity to mitigate the downstream instability caused by high KL divergence. As demonstrated in Figure 4 and Table 2, GEPO's group expectation mechanism effectively suppresses the explosion of importance weight variance even when KL divergence is high. This allows GEPO to maintain stable training and avoid collapse in many scenarios where baseline methods (like GRPO and GSPO) fail — thereby establishing algorithmic robustness to policy divergence as a core contribution of this work.

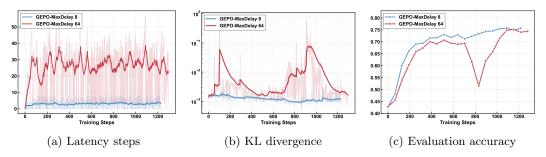


Figure 7: Training processes under different latency conditions

5 Related Work

In recent years, reinforcement learning has become central to post-training LLMs (Ziegler et al., 2019). Researchers have identified efficiency bottlenecks in traditional synchronous RL frameworks. IMPALA(Espeholt et al., 2018) is a centralized, high-throughput distributed reinforcement learning architecture designed for traditional RL settings—not LLMs—that uses V-trace to correct policy lag between asynchronous actors and a central learner. Wu et al. (2025) first theoretically explored asynchronous RLHF (Ouyang et al., 2022), proposing to decouple generation and training across GPU clusters, enabling scalable and efficient RL fine-tuning of large language models up to 405B parameters with provable speedup over synchronous baselines. However, it focuses on single-machine or small clusters, neglecting dynamic network delays in heterogeneous environments. To address practical efficiency, AREAL (Fu et al., 2025) fully decouples generation and training, using staleness thresholds and a decoupled PPO objective (Schulman et al., 2017b) to handle outdated samples. It improves both training speed and final performance on reasoning and code tasks—but assumes stable networks, unlike the unpredictable, high-latency settings targeted by HeteroRL. Prime Intellect (Team et al., 2025) offers a decentralized, asynchronous framework for community compute, ensuring trust via verifiable inference, stability via two-sided GRPO clipping, and controllable reasoning with length-aware rewards. These works motivate our design: a robust, delay-tolerant framework for heterogeneous, geographically distributed RL.

6 Conclusion

We propose HeteroRL, a heterogeneous reinforcement learning framework designed for training LLMs across geographically distributed and resource-heterogeneous nodes, paired with GEPO—a novel policy optimization algorithm that stabilizes training under high latency. By decoupling rollout sampling from parameter updates, HeteroRL eliminates synchronization bottlenecks inherent in traditional RL pipelines. GEPO addresses the explosion of variance of importance weight caused by stale policies through group expectation importance weight, provably reducing variance exponentially, particularly under large KL divergence between the sampling and learning policies. This work establishes a practical foundation for scalable, delay-tolerant decentralized RL, making it well-suited for real-world LLM post-training in heterogeneous, wide-area network environments.

ETHICS STATEMENT

This work presents a systems and algorithmic contribution to stable reinforcement learning in heterogeneous, geographically distributed environments. It uses only public, non-sensitive datasets (e.g., MATH, AIME) and does not involve human subjects, personal data, or high-risk applications. The proposed method (GEPO) is task-agnostic and neutral in intent; it does not introduce bias or discrimination, as training focuses solely on mathematical reasoning.

REPRODUCIBILITY STATEMENT

We provide full implementation details in Appendix B.1, including model configurations, training hyperparameters, and latency simulation setup. All datasets used (MATH, AMC, AIME) are public. The algorithm (GEPO) is precisely described in Section 3 and Listing 1, with theoretical guarantees in Appendix A. Anonymous code and configuration files are included in the supplementary materials.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- AIME. Aime 2024. https://huggingface.co/datasets/HuggingFaceH4/aime_2024, 2024. Accessed: 2025-03-17.
- AIME. Aime 2025. https://huggingface.co/datasets/yentinglin/aime_2025, 2025. Accessed: 2025-03-17.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *ArXiv Preprint ArXiv:2212.08073*, 2022.
- Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. Safe RLHF: Safe reinforcement learning from human feedback. In *Proceedings of the Twelfth International Conference on Learning Representations*, 2024.
- Muzhi Dai, Shixuan Liu, and Qingyi Si. Stable reinforcement learning for efficient reasoning, 2025. URL https://arxiv.org/abs/2505.18086.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *ArXiv Preprint ArXiv:2407.21783*, 2024.
- Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Volodymir Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, Shane Legg, and Koray Kavukcuoglu. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures, 2018. URL https://arxiv.org/abs/1802.01561.
- Wei Fu, Jiaxuan Gao, Xujie Shen, Chen Zhu, Zhiyu Mei, Chuyi He, Shusheng Xu, Guo Wei, Jun Mei, Jiashu Wang, et al. Areal: A large-scale asynchronous reinforcement learning system for language reasoning. arXiv preprint arXiv:2505.24298, 2025.
- Daya Guo, Dejian Yang, Haowei Zhang, et al. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081):633–638, 2025. doi: 10.1038/s41586-025-09422-z. URL https://doi.org/10.1038/s41586-025-09422-z.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. arXiv preprint arXiv:2103.03874, 2021.

- Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif Rasul, Longhui Yu, Albert Q Jiang, Ziju Shen, et al. Numinamath: The largest public dataset in ai4maths with 860k pairs of competition math problems and solutions. Hugging Face repository, 13(9):9, 2024.
 - Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective. arXiv preprint arXiv:2503.20783, 2025.
 - Michael Noukhovitch, Shengyi Huang, Sophie Xhonneux, Arian Hosseini, Rishabh Agarwal, and Aaron Courville. Asynchronous rlhf: Faster and more efficient off-policy rl for language models. arXiv preprint arXiv:2410.18252, 2024.
 - Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
 - John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *ArXiv Preprint ArXiv:1707.06347*, 2017a.
 - John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017b.
 - Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. arXiv preprint arXiv:2402.03300, 2024.
 - Vaishnavi Shrivastava, Ahmed Awadallah, Vidhisha Balachandran, Shivam Garg, Harkirat Behl, and Dimitris Papailiopoulos. Sample more to think less: Group filtered policy optimization for concise reasoning. arXiv preprint arXiv:2508.09726, 2025.
 - Ziang Song, Tianle Cai, Jason D. Lee, and Weijie J. Su. Reward collapse in aligning large language models, 2023. URL https://arxiv.org/abs/2305.17608.
 - Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback. In *Proceedings of the International Conference on Neural Information Processing Systems*, pp. 3008–3021, 2020.
 - Prime Intellect Team, Sami Jaghouar, Justus Mattern, Jack Min Ong, Jannik Straube, Manveer Basra, Aaron Pazdera, Kushal Thaman, Matthew Di Ferrante, Felix Gabriel, et al. Intellect-2: A reasoning model trained through globally decentralized reinforcement learning. arXiv preprint arXiv:2505.07291, 2025.
 - Bo Wu, Sid Wang, Yunhao Tang, Jia Ding, Eryk Helenowski, Liang Tan, Tengyu Xu, Tushar Gowda, Zhengxing Chen, Chen Zhu, et al. Llamarl: A distributed asynchronous reinforcement learning framework for efficient large-scale llm trainin. arXiv preprint arXiv:2505.24034, 2025.
 - Changyi Xiao, Mengdi Zhang, and Yixin Cao. Bnpo: Beta normalization policy optimization. arXiv preprint arXiv:2506.02864, 2025.
 - An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. arXiv preprint arXiv:2505.09388, 2025.

- Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the wild. arXiv preprint arXiv:2503.18892, 2025.
- Han Zhang, Yu Lei, Lin Gui, Min Yang, Yulan He, Hui Wang, and Ruifeng Xu. Cppo: Continual learning for reinforcement learning with human feedback. In *The Twelfth International Conference on Learning Representations*, 2024.
- Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, et al. Group sequence policy optimization. arXiv preprint arXiv:2507.18071, 2025.
- Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. arXiv preprint arXiv:1909.08593, 2019.

A THEORETICAL PROOF OF IMPORTANCE SAMPLING VARIANCE

This appendix analyzes a newly proposed importance sampling weight $w_{\text{new}}(x) = \frac{p(x)}{\mathbb{E}_q[q]}$, where $\mathbb{E}_q[q] = \int q(x)^2 dx$, and compares its variance with the standard importance sampling weight $w_{\text{std}}(x) = \frac{p(x)}{q(x)}$.

Problem Setting Let the target distribution be p(x) and the proposal distribution be q(x). We aim to estimate:

$$\mathbb{E}_p[f] = \int f(x)p(x)dx. \tag{6}$$

Since direct sampling from p is difficult, we employ importance sampling by drawing samples from q.

Standard Importance Sampling The standard weight is defined as:

$$w_{\rm std}(x) = \frac{p(x)}{q(x)}. (7)$$

Its expectation under q is:

$$\mathbb{E}_q\left[\frac{p(x)}{q(x)}\right] = \int \frac{p(x)}{q(x)} q(x) dx = \int p(x) dx = 1,\tag{8}$$

thus it is unbiased. Its variance is:

$$\operatorname{Var}_{q}(w_{\mathrm{std}}) = \mathbb{E}_{q} \left[\left(\frac{p}{q} \right)^{2} \right] - \left(\mathbb{E}_{q} \left[\frac{p}{q} \right] \right)^{2}$$

$$= \int \frac{p(x)^{2}}{q(x)} dx - 1.$$
(9)

Denoted as:

$$Var_{std} = \int \frac{p(x)^2}{q(x)} dx - 1.$$
 (10)

Group Expectation Importance Sampling The new weight is defined as:

$$w_{\text{new}}(x) = \frac{p(x)}{\mathbb{E}_q[q]}, \quad \text{where} \quad \mathbb{E}_q[q] = \int q(x)^2 dx.$$
 (11)

Its expectation is:

$$\mathbb{E}_q[w_{\text{new}}] = \frac{1}{\mathbb{E}_q[q]} \int p(x)q(x)dx = \frac{\langle p, q \rangle}{\|q\|_2^2},\tag{12}$$

where $\langle p,q \rangle$ denotes the inner product $\int p(x)q(x)dx$. Generally, $\langle p,q \rangle \neq ||q||_2^2$, making this estimator biased. Its variance is:

$$\operatorname{Var}_{q}(w_{\text{new}}) = \mathbb{E}_{q} \left[\left(\frac{p(x)}{\mathbb{E}_{q}[q]} \right)^{2} \right] - \left(\mathbb{E}_{q} \left[\frac{p(x)}{\mathbb{E}_{q}[q]} \right] \right)^{2}$$

$$= \frac{1}{(\mathbb{E}_{q}[q])^{2}} \left(\int p(x)^{2} q(x) dx - \left(\int p(x) q(x) dx \right)^{2} \right).$$
(13)

Denoted as:

$$Var_{new} = \frac{1}{\left(\int q(x)^2 dx\right)^2} \left(\int p(x)^2 q(x) dx - \left(\int p(x) q(x) dx\right)^2\right). \tag{14}$$

Variance Comparison We compare:

$$\operatorname{Var}_{\text{std}} = \int \frac{p(x)^2}{q(x)} dx - 1,$$

$$\operatorname{Var}_{\text{new}} = \frac{1}{\left(\int q(x)^2 dx\right)^2} \left(\int p(x)^2 q(x) dx - \left(\int p(x) q(x) dx\right)^2\right).$$
(15)

A.1 Variance Comparison in Discrete Space

Since the action space of large models is discrete, this section discusses the variance difference $\Delta = Var_{std} - Var_{new}$ in discrete probability space. The integral expressions in the continuous case naturally transition to discrete summation forms:

- Replace continuous integrals $\int dx$ with discrete summations $\sum_{i=1}^{n}$;
- Replace probability density functions p(x), q(x) with probability masses p_i , q_i ;
- Maintain the structural form of variance expressions unchanged.

Notation and Setting Let the sample space be a finite set $\mathcal{X} = \{1, 2, \dots, n\}$, where $n \geq 2$. Let $p = (p_1, \ldots, p_n)$ and $q = (q_1, \ldots, q_n)$ be two probability distributions satisfying:

• $p_i > 0$, $\sum_{i=1}^n p_i = 1$,

 • $q_i > 0$, $\sum_{i=1}^n q_i = 1$.

Define the following four key quantities, corresponding to the integral terms in the continuous expressions:

$$A = \left(\sum_{i=1}^{n} q_i^2\right)^2, \qquad \text{(corresponding to } \left(\int q(x)^2 dx\right)^2) \qquad (16)$$

$$B = \left(\sum_{i=1}^{n} p_i q_i\right)^2, \qquad \text{(corresponding to } \left(\int p(x) q(x) dx\right)^2) \qquad (17)$$

$$I_1 = \sum_{i=1}^n \frac{p_i^2}{q_i}, \qquad (corresponding to \int \frac{p(x)^2}{q(x)} dx) \qquad (18)$$

$$I_2 = \sum_{i=1}^n p_i^2 q_i, \qquad \text{(corresponding to } \int p(x)^2 q(x) dx)$$
 (19)

Accordingly, the variance difference can be written as:

$$\Delta = I_1 + \frac{B - A - I_2}{A}.\tag{20}$$

Lemma 1 (Range of Quantities). Under the above setting, we have:

- $A \in \left[\frac{1}{n^2}, 1\right]$,
- $B \in [0, 1],$
- $I_1 \in [1, \infty)$,
- $I_2 \in (0, 1]$.

Proof.

- **1. Range of** A: By the power mean inequality, $\sum q_i^2 \ge \frac{1}{n}$, and $\sum q_i^2 \le 1$, thus $A \in [1/n^2, 1]$.
- **2.** Range of $B: \sum_{i=1}^{\infty} p_i q_i \in (0,1]$, thus $B \in [0,1]$. **3.** Range of I_1 : By the Cauchy-Schwarz inequality, $I_1 \geq 1$, and it can approach infinity.
- **4. Range of** I_2 : Since $p_i^2 \le p_i$, we have $I_2 \le 1$, and $I_2 > 0$.

Theorem 1. Let p,q be discrete probability distributions. Then there exists a constant Csuch that:

$$\operatorname{Var}\left[\frac{p(y|x)}{q(y|x)}\right] - \operatorname{Var}\left[\frac{p(y|x)}{\widehat{\mathbb{E}}_{q}[q(y|x)]}\right] \ge \exp\left(D_{\mathrm{KL}}(p\|q)\right) - C. \tag{21}$$

In particular, when $D_{\mathrm{KL}}(p||q) > \log C$, it holds that $\mathrm{Var}[\frac{p(y|x)}{q(y|x)}] > \mathrm{Var}[\frac{p(y|x)}{\widehat{\mathbb{R}}_{-}[q(y|x)]}]$.

Proof. **Step 1.** From the fundamental inequality relationship between KL divergence and χ^2 divergence (Pinsker's inequality):

$$D_{\text{KL}}(p||q) \le \log(1 + D_{\chi^2}(p||q)),$$
 (22)

where the chi-square divergence is defined as:

$$D_{\chi^2}(p||q) = \sum_{i=1}^n \frac{(p_i - q_i)^2}{q_i} = \sum_{i=1}^n \frac{p_i^2}{q_i} - 1 = I_1 - 1.$$
 (23)

Substituting yields:

$$D_{\mathrm{KL}}(p||q) \le \log(I_1),\tag{24}$$

therefore:

$$I_1 \ge \exp\left(D_{\mathrm{KL}}(p\|q)\right). \tag{25}$$

Step 2. From Lemma 1, we know that A, B, and I_2 satisfy the following bounds:

$$A \in \left[\frac{1}{n^2}, 1\right],\tag{26}$$

$$B \in [0, 1], \tag{27}$$

$$I_2 \in (0,1].$$
 (28)

Consider the lower bound of the expression $\frac{B-A-I_2}{A}$. To obtain its minimum value, we take:

- Minimum value of B: B = 0
- Minimum value of A: $A = \frac{1}{n^2}$
- Maximum value of I_2 : $I_2 = 1$

Substituting yields:

$$\frac{B-A-I_2}{A} \ge \frac{0-1/n^2-1}{1/n^2} = -(n^2+1). \tag{29}$$

Step 3. Substituting inequalities (1) and (2) into the expression for Δ :

$$\Delta = I_1 + \frac{B - A - I_2}{A} \ge \exp\left(D_{KL}(p||q)\right) - (n^2 + 1). \tag{30}$$

When $D_{KL}(p||q) > \log(n^2 + 1)$, we have:

$$\exp(D_{KL}(p||q)) > n^2 + 1,$$
 (31)

thus:

$$\Delta > 0, \tag{32}$$

i.e., $Var_{std} > Var_{new}$.

Corollary 1. In discrete space, if $D_{KL}(p||q) > \log(n^2 + 1)$, then:

$$Var_{std} > Var_{new},$$
 (33)

i.e., the new estimator has smaller variance.

It should be noted that A only attains the value $\frac{1}{n^2}$ when q follows a uniform distribution. In practice, when large models generate responses, the distribution tends to be long-tailed, so the value of A is much greater than $\frac{1}{n^2}$, and the constant $C_{real} \ll \log(n^2 + 1)$. For example, we randomly generated 128 tokens using Qwen3-1.7B (n = 151936), and the standard variance and average of A was $0.432_{\pm 0.36} \gg \frac{1}{n^2}$.

B Supplementary experiments

B.1 Implementation Details

All experiments use the Qwen3-1.7B/8B model with a maximum input length of 768 and output length of 2048/4096 tokens under both think and no-think mode, limited by computational constraints and low token efficiency (Shrivastava et al., 2025) (reward/length) at full context length¹. Training follows a GRPO-like algorithm with a learning rate of 1×10^{-6} , 3% linear warmup, per-device batch size 8, and gradient accumulation of 1, with gradient checkpointing enabled for memory efficiency. Evaluations occur every 32 or 64 steps. To model network latency in heterogeneous environments, we introduce a log-normal delay simulator bounded between 60 and 1800 seconds (99.5% CI), with default delay at 60 seconds and policy staleness varied across 0–64 effective steps. For online training, KL divergence is not used, in order to better evaluate the training stability of the algorithms. For heterogeneous settings, CPPO-KL (Zhang et al., 2024) loss with coefficient 0.005 is applied. Another reason for using CPPO-KL is memory efficiency, as it does not require a separate reference model. Rollouts are generated using vLLM with 8 parallel responses per prompt, and each run lasts 3 epochs, with metrics logged via Weights & Biases. The system prompt used in the experiments is shown in Figure 8.

[think mode] You are a helpful AI Assistant, designed to provided well-reasoned and detailed responses. You FIRST think about the reasoning process as an internal monologue and then provide the user with the answer. Please put your final answer within \boxed{}. Also, indicate that it is the answer.

[no-think mode] You are a helpful AI Assistant, designed to provided well-reasoned and detailed responses. Please put your final answer within \boxed{}. Also, indicate that it is the answer.

Figure 8: System prompt of all trainings in our experiments.

B.2 Baselines

The baseline methods compared in our experiments are as follows:

- Group Relative Policy Optimization (GRPO) (Shao et al., 2024) is a reinforcement learning algorithm that enhances mathematical reasoning in LLMs by estimating advantages through group-relative reward normalization—comparing responses within a group to the same query—thereby eliminating the need for a separate value network and reducing memory overhead compared to PPO.
- Dr. GRPO (Liu et al., 2025) is a debiased variant of GRPO that removes the per-response length normalization and per-question reward standard deviation normalization, thereby eliminating optimization biases that artificially inflate response length and improving token efficiency while preserving reasoning performance.
- Beta Normalization Policy Optimization (BNPO) (Xiao et al., 2025) is a novel reinforcement learning algorithm that dynamically normalizes binary rewards using an adaptively parameterized Beta distribution to reduce gradient variance and enhance training stability for large language models.
- Group Sequence Policy Optimization (GSPO) (Zheng et al., 2025) is a novel reinforcement learning algorithm for large language models that defines importance
 ratios based on sequence likelihood and performs sequence-level clipping and optimization, thereby achieving superior training stability, efficiency, and performance
 compared to token-level methods like GRPO.

¹This 2k/4k-token limit balances cost and efficiency. Longer outputs increase memory and training time, making high-latency experiments impractical. Crucially, as shown in Figure 9(c), "think mode" yields a high "Overlength Ratio"—most long outputs are truncated and wasted. Recent studies (Shrivastava et al., 2025; Dai et al., 2025) consistently show that reasoning does not require excessively long chains of thought; redundant thinking merely wastes resources. Thus, 2048 tokens ensure fair, stable, and manageable experiments.

B.3 Supplementary Results of Main Experiments

As shown in Table 3, GEPO outperforms all baselines—including GRPO and GSPO—in both best and final performance across zero-delay and high-delay settings, demonstrating superior effectiveness and training stability. A critical observation from the results is that, although GSPO's technical report claims improved stability, we only observe GSPO to be more stable than GRPO under the Online RL setting.

Table 3: Performance of GEPO and baseline methods under Online RL and Hetero RL scenarios (2k limiation).

| Method | AMC | AMC2023 | | AIME2024 | | AIME2025 | | MATH500 | | rage |
|------------------------------------|------|---------|------|----------|------|----------|------|-------------|------|------|
| 1,1001101 | Best | Last | Best | Last | Best | Last | Best | Last | Best | Last |
| Qwen3-1.7B | 25.6 | - | 1.6 | - | 3.9 | - | 54.7 | - | 21.5 | - |
| Max Tolerable Delay 0 (Online RL) | | | | | | | | | | |
| BNPO | 54.3 | 0.0 | 18.4 | 0.0 | 19.1 | 0.0 | 78.7 | 0.0 | 42.6 | 0.0 |
| Dr.GRPO | 53.4 | 14.3 | 19.1 | 1.6 | 18.8 | 2.0 | 78.6 | 35.9 | 42.5 | 13.5 |
| GRPO | 56.3 | 23.4 | 20.7 | 0.4 | 19.9 | 2.3 | 79.8 | 49.7 | 44.2 | 19.0 |
| GSPO | 54.1 | 27.8 | 23.8 | 3.1 | 20.7 | 4.3 | 79.9 | 62.1 | 44.6 | 24.3 |
| GEPO (ours) | 56.9 | 56.9 | 21.9 | 16.4 | 20.3 | 14.1 | 80.4 | 78.1 | 44.9 | 41.4 |
| Max Tolerable Delay 64 (Hetero RL) | | | | | | | | | | |
| BNPO | 45.0 | 43.1 | 12.1 | 11.3 | 12.5 | 10.1 | 71.1 | 69.3 | 35.2 | 33.5 |
| Dr.GRPO | 48.4 | 48.4 | 17.2 | 17.2 | 14.8 | 14.8 | 73.9 | 73.9 | 38.6 | 38.6 |
| GRPO | 46.6 | 46.6 | 19.1 | 14.5 | 14.8 | 14.8 | 74.9 | 74.9 | 38.9 | 37.7 |
| GSPO | 54.4 | 23.8 | 17.6 | 1.6 | 17.6 | 2.7 | 78.2 | 55.6 | 42.0 | 20.9 |
| GEPO (ours) | 53.8 | 53.8 | 21.9 | 21.9 | 18.8 | 18.8 | 79.6 | 79.6 | 43.5 | 43.5 |

B.4 Ablation Study

Since the removal of the variance divisor term in the advantage function has already been extensively validated in prior work (Liu et al., 2025), we focus solely on ablating the Group Expectation component. Table 4 compares three importance weights (IW) in Listing 1, namely token-level (IW of GRPO), sequence-level (IW of GSPO), and group-level (IW of GEPO), under a max tolerable delay of 64 steps. Table 4 presents the results, the sequence-level IW proposed by GSPO does not bring significant stability improvement. Although it outperforms token-level weighting in terms of best performance, it suffers from severe performance degradation by the end of training.

Table 4: Ablation study of different importance weights (Hetero RL mode).

| Ablation | AMC2023 | | AIME2024 | | AIME2025 | | MATH500 | | Average | |
|--------------------|--------------|--------------|----------|------|--------------|------|---------|------|--------------|--------------|
| 1151401011 | Best | Last | Best | Last | Best | Last | Best | Last | Best | Last |
| group-lv | 53.8 | 53.8 | 21.9 | 21.9 | 18.8 | 18.8 | 79.6 | 79.6 | 43.5 | 43.5 |
| token-lv seq-lv | 46.1 55.2 | 43.9 24.1 | | | 15.3 18.1 | - | | | 38.6 42.0 | 36.8 21.3 |

B.5 Comparison of Think and Non-think Mode

In *non-think mode*, all methods struggle with exploration, shown by consistently lower last scores vs. best — indicating late-stage instability.

Notably:

- GEPO (ours) dominates across all benchmarks and delays in both best and last, proving superior efficiency and stability even without thinking steps.
- At zero delay, GEPO beats top baseline (GSPO) by +3.7 (best) and +20.8 (last), resisting early collapse.

- With delay=64, baselines improve (e.g., Dr.GRPO: $0.0 \rightarrow 33.4$), showing slack aids stability. GEPO still leads (last = 38.0 vs. 35.7).
- Vanilla Qwen3-1.7B (best=33.2) is outperformed by all RL methods confirming RL's value even without reasoning.

In short: while non-think mode limits reasoning, GEPO delivers unmatched stability and final performance, ideal for latency-sensitive or real-time settings.

Table 5: Performance comparison using Qwen3-1.7B (non-think, 2k).

| Method | AMC2023 | | AIME2024 | | AIME2025 | | MATH500 | | Average | |
|------------------------|----------------------------|------|----------------------------|------|----------------------------|------|----------------------------|------|----------------------------|------|
| | $\overline{\mathrm{Best}}$ | Last |
| Qwen3-1.7B | 42.8 | - | 10.2 | - | 9.4 | - | 70.2 | - | 33.2 | - |
| Max Tolerable Delay 0 | | | | | | | | | | |
| BNPO | 43.7 | 0.0 | 13.7 | 0.0 | 13.3 | 0.0 | 74.2 | 0.4 | 36.2 | 0.1 |
| Dr.GRPO | 45.0 | 0.0 | 14.4 | 0.0 | 11.3 | 0.0 | 73.6 | 0.0 | 36.1 | 0.0 |
| GRPO | 50.0 | 28.8 | 16.4 | 7.8 | 13.7 | 7.2 | 77.5 | 59.7 | 39.4 | 25.9 |
| GSPO | 53.1 | 23.8 | 14.4 | 2.7 | 16.0 | 0.4 | 76.5 | 55.6 | 40.0 | 20.6 |
| GEPO (ours) | 55.0 | 52.5 | 22.3 | 22.3 | 18.4 | 13.7 | 79.2 | 77.1 | 43.7 | 41.4 |
| Max Tolerable Delay 64 | | | | | | | | | | |
| BNPO | 42.1 | 28.6 | 11.1 | 4.7 | 10.2 | 6.4 | 68.3 | 36.1 | 32.9 | 19.0 |
| Dr.GRPO | 45.0 | 41.6 | 14.1 | 9.7 | 13.3 | 10.1 | 72.5 | 72.3 | 36.2 | 33.4 |
| GRPO | 46.3 | 46.3 | 14.5 | 14.5 | 13.6 | 10.1 | 72.6 | 71.7 | 36.8 | 35.7 |
| GSPO | 47.2 | 36.6 | 13.7 | 4.3 | 13.3 | 8.9 | 75.1 | 67.9 | 37.3 | 29.4 |
| GEPO (ours) | 52.5 | 52.5 | 14.5 | 10.2 | 14.5 | 12.1 | 77.4 | 77.1 | 39.7 | 38.0 |

A comparison of RL under think vs. non-think modes reveals key trade-offs in reasoning, stability, cost, and performance.

- 1) Performance and Stability (Delay = 0) GEPO in non-think achieves Best = 43.7 and Last = 41.4 — near-perfect retention. In think mode, it reaches Best = 44.9 (+1.2) but identical Last = 41.4 — no final gain. Baselines like GRPO collapse $(44.2 \rightarrow$ 19.0), showing thinking destabilizes training without proper control.
- 2) With Delay = 64 In non-think, methods rebound sharply: Dr.GRPO $(0.0 \rightarrow 33.4)$, GRPO $(25.9 \rightarrow 35.7)$ — minimal slack prevents collapse. In think, only GEPO retains peak perfectly (Last = Best = 43.5); others degrade (e.g., GSPO: $41.9 \rightarrow 20.9$). \Rightarrow Thinking introduces instability unless explicitly regularized.
- 3) Exploration vs. Efficiency (Fig. 9) Think produces longer rollouts deeper reasoning — but with higher overlength ratios, risking wasted compute and divergence. non-think yields shorter, efficient trajectories with lower overlength — yet matches or exceeds performance when paired with stable optimizers like GEPO.

Overall, think offers marginal peak gains at the cost of instability and overhead. Non-think + GEPO is more stable, efficient, and often equally effective — ideal for real-time or latency-sensitive deployment. Choose based on delay tolerance, interpretability needs, and required training stability. Figure 9 compares rollout lengths in GEPO under think vs. non-think modes, with three subplots:

- 1) Average Length: The red curve (think) stays consistently above blue (non-think), showing longer rollouts — likely due to added reasoning or policy deliberation.
- 2) Terminated Length: think again produces significantly longer successful rollouts, suggesting it enables deeper, valid exploration.
- 3) Overlength Ratio: think incurs a much higher overlength rate, especially early in training — indicating that while it boosts exploration, it also risks inefficient, overly long sequences.

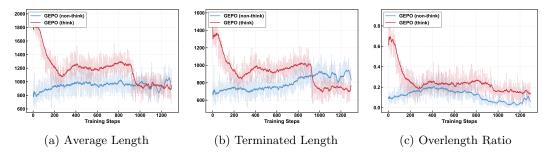


Figure 9: Comparison of rollout (cliped by 2k tokens) lengths under think and non-think modes.

Takeaway: think extends useful exploration and improves completion quality, but at the cost of higher computational overhead. In practice, this trade-off can be managed via length limits or dynamic rollout control to preserve efficiency without sacrificing performance.

C NETWORK LATENCY

We test the performance limits of GEPO and baseline methods in simulated network delay scenarios and train the model on heterogeneous computing resources connected through a real-world internet connection.

C.1 SIMULATED NETWORK LATENCY CONFIGURATION

Our goal is to simulate RL training of large models over internet-connected heterogeneous compute clusters, where network delays are inherently uncertain. To evaluate algorithmic performance and training stability under extreme and variable latency conditions, we employ three widely used delay distributions P_d : log-normal, Weibull, and exponential. We set a high maximum delay threshold of 1800 seconds, which is sufficient to cover typical model and data transmission times. In our simulation, model parameter synchronization and rollout data transfer are implemented as follows:

- Learner saves model: The learner periodically saves model checkpoints to a shared model file path Model_Sync_Path via torch.save_pretrained().
- Sampler loads model: The sampler generates data using the current model until a delay D_M , drawn from P_d , elapses. It then loads the latest model from Model_Sync_Path. The sampler remains active throughout—no idling occurs.
- Sampler sends data: Generated rollout batches are saved to Rollout_Sync_Path with timestamp T_{sync} . To simplify implementation, data transmission is assumed instantaneous; its latency is effectively merged into the model sync delay D_M , without affecting simulation validity.
- Learner uses data: The learner trains on rollout data from Rollout_Sync_Path that falls within a recent time window (e.g., no older than 1800 seconds).

C.2 Real-World Network Scenarios

To evaluate the performance of heterogeneous reinforcement learning algorithms in realistic network environments, we develop a communication toolkit based on ZeroMQ, supporting TCP/IP-based transmission of inference trajectories from samplers to learners and synchronized model parameter updates. As shown in Figure 10, the toolkit enables multi-node communication over wide-area networks (WANs), with the following core communication logic:

 The learner continuously listens for and buffers inference trajectory messages from samplers. It automatically updates its trajectory buffer upon message arrival and broadcasts the latest model parameters to all connected samplers once a predefined synchronization interval is reached.

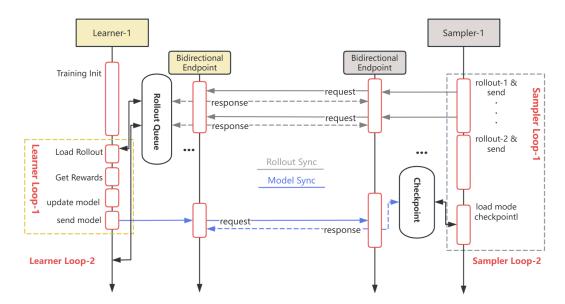


Figure 10: Communication flow between a learner and sampler in a distributed reinforcement learning system, showing trajectory collection and model synchronization over a bidirectional endpoint.

 After generating an inference trajectory, the sampler sends it to the learner and continuously listens for parameter update messages. Upon receiving new model parameters, the sampler updates its local parameters and resumes sampling before the next sampling step.

Key features of the communication toolkit include:

- 1) Many-to-Many Communication Pattern The toolkit supports elastic node scaling with dynamically reconfigurable routing topologies, enabling seamless adaptation to node join/leave events. This facilitates efficient distributed parameter synchronization and trajectory collection in large-scale deployments.
- 2) Chunked Message Transmission Inspired by Shardcast, the toolkit employs adaptive message chunking for both model parameters and trajectory data. The chunk size is dynamically adjusted according to real-time bandwidth conditions. At the receiver side, chunks are reassembled and integrity-verified, effectively mitigating the impact of high latency in wide-area networks.
- **3)** Communication Safety Mechanisms Thread-safe message passing across multiple processes is ensured via a double-buffering queue design and fine-grained locking. This prevents race conditions during concurrent read/write operations and guarantees reliable parameter synchronization.
- 4) Configurable Routing Topology The toolkit supports user-defined connectivity patterns between learners and samplers, accommodating asymmetric node configurations (e.g., one learner connected to four samplers). This enhances communication flexibility and scalability in dynamic, heterogeneous environments.

D Engineering Optimization: Localized Reward Computation for Reduced Communication

In large-scale heterogeneous reinforcement learning (HeteroRL), where sampler and learner nodes are geographically distributed, network communication overhead becomes a critical bottleneck. A major source of this overhead lies in the reward aggregation phase, where traditional implementations require an all_gather operation across all processes to collect rewards for group-wise normalization (e.g., computing mean and standard deviation per

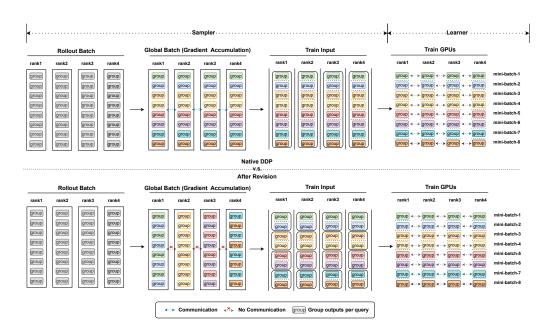


Figure 11: Native Distributed Data Parallel and Revision.

group of generations). This global synchronization introduces significant latency, especially under high network delay or when using large numbers of distributed nodes. To address this, we introduce a key engineering optimization: **Localized Reward Computation**. Instead of gathering rewards from all processes, we ensure that each group of generations (e.g., *G* responses per prompt) is entirely generated and scored *within the same process or node*. This design guarantees that all samples belonging to a single group reside locally, enabling the calculation of group statistics (mean, std) without any cross-process communication.

This optimization is particularly synergistic with GEPO's design philosophy. GEPO mitigates training instability under high latency by reducing the variance of importance weights through group-level expectation. Our engineering improvement complements this by reducing the system's sensitivity to communication latency itself, creating a virtuous cycle: the algorithm is robust to policy staleness, and the system is optimized to minimize the communication that causes staleness.

The core change is implemented by removing the global gather operation in the reward calculation function. This seemingly minor change yields substantial performance benefits. The following table summarizes the key improvements:

Table 6: Comparison of Communication Overhead: Before vs. Optimized (Localized Reward)

| Component | Before Optimization | After Optimization (Ours) | | |
|-------------------------|-------------------------------------|---------------------------------|--|--|
| Reward Aggregation | Requires all_gather for every batch | No communication required | | |
| Group Statistics | Computed globally across all GPUs | Computed locally per GPU | | |
| Communication Frequency | High (per batch) | None for reward calculation | | |
| Latency Sensitivity | High (blocked by slowest node) | Low (fully asynchronous) | | |
| Scalability | Limited by network bandwidth | Highly scalable with node count | | |

By eliminating this frequent and costly synchronization point, our system achieves higher throughput and better resource utilization. This is especially crucial in the HeteroRL setting, where network conditions are unpredictable. The optimization allows sampler nodes to

operate more independently, reducing their idle time waiting for global synchronization and making the overall training pipeline more resilient to network fluctuations. This engineering refinement, combined with the GEPO algorithm, provides a holistic solution for stable and efficient training in truly heterogeneous, high-latency environments.

E From Token-Level to Group-Level Importance Weight

Token-Level Importance Weight In traditional policy optimization methods, such as GRPO (Shao et al., 2024) or PPO (Schulman et al., 2017a), importance sampling is typically performed at the token-level. Specifically, for a generated sequence $y = (y_1, \ldots, y_T)$, the importance weight is computed token by token:

$$w_{token} = \operatorname{clip}\left(\frac{p(y_t \mid x, y_{< t})}{q(y_t \mid x, y_{< t})}, 1 - \epsilon, 1 + \epsilon\right), \tag{34}$$

and used to compute the policy gradient at each time step. However, this per-token reweighting scheme suffers from two key limitations:

- 1) Inconsistency between optimization and reward granularity: Concurrent work such as GSPO (Zheng et al., 2025) argues that since the reward is assigned at the sequence level—i.e., to the full response y—the importance weighting for policy updates should also operate at the same granularity, rather than at the token level.
- 2) High variance in token-level probabilities: Beyond this alignment issue, we provide an explanation from the perspective of importance weight variance: because the final reward depends on the entire sample, large local changes in token probabilities can lead to extreme values in the ratio $\frac{p(y_t|x,y_{< t})}{q(y_t|x,y_{< t})}$. This inflates the variance of token-level importance weights, causing them to frequently fall outside the clipping range. Once clipped, gradients for these tokens are effectively zeroed out due to the stop-gradient behavior of torch.clamp, preventing meaningful updates. As a result, tokens that require large corrections may be ignored, leading to inefficient optimization.

Sample-Level Importance Weight The sample-level importance weighting treats the entire response as a sampling unit and computing the weight based on its full conditional probability. Given a prompt x, both the target policy p(y|x) and the behavior policy q(y|x) are defined as the joint probability of the sequence:

$$p(y|x) = \left| \prod_{t=1}^{T} p(y_t \mid x, y_{< t}) \right|^{\frac{1}{T}}, \quad q(y|x) = \left| \prod_{t=1}^{T} q(y_t \mid x, y_{< t}) \right|^{\frac{1}{T}}.$$
 (35)

The sample-level importance weight is then:

$$w_{\text{sequence}}(y|x) = \text{clip}\left(\frac{p(y|x)}{q(y|x)}, 1 - \epsilon, 1 + \epsilon\right). \tag{36}$$

This formulation computes the full sequence-level ratio before clipping, thereby avoiding premature truncation caused by high-variance individual tokens. By preserving gradient flow across the entire sequence, sample-level weighting enables more stable and effective policy updates, especially under high policy divergence induced by network latency.

Group-Level Importance Weight Group-level importance weighting represents a paradigm shift in reinforcement learning optimization by recognizing that policy updates should not only align with the granularity of reward assignment but also leverage higher-order statistical relationships among multiple samples. The key insight is that individual samples should not be treated in isolation; instead, their collective behavior under the same prompt provides crucial information for stable policy updates. By considering the expected value of proposal probabilities within a group of responses, we effectively smooth out the erratic fluctuations that plague token-level and even sequence-level weighting schemes. This approach acknowledges an important reality of distributed training: in heterogeneous environments with network latency, policy divergence is inevitable, and our optimization methods must be designed to gracefully handle—not merely tolerate—this divergence. The group-level perspective transforms what was previously seen as a limitation (policy staleness due to latency) into an opportunity for more robust learning through statistical regularization of importance weights.

Table 7: Core Characteristics of Online, Offline, and Heterogeneous Reinforcement Learning Paradigms

| Aspect | Online RL | Offline RL | Heterogeneous RL |
|---|---|--|---|
| Data Generation | Online interaction: Data generated instantly by the current policy . | Fixed, pre-collected dataset: Uses static data collected by some (unknown) behavior policy. | Delayed interaction: Data generated by historical policy versions (due to unpredictable network delay). |
| Policy Version during Data Generation | Always current: Requires strict synchronization with the learning policy. | Fixed: The behavior policy is fixed and inherent to the dataset. | Dynamically stale: Staleness determined by unpredictable network delay (core characteristic). |
| System Architecture | Tightly-coupled: Actor (environment interaction) and Learner (parameter update) typically co-located or on a low-latency network. | Single-machine or simple distributed: No real-time interaction required; training is data-driven. | Geographically decoupled: Actor and Learner separated by a high-latency network, tolerant to delays. |
| Handling of System Latency | Treated as failure: Requires synchronization; latency causes resource idling and training stalls. | Not applicable: Training process has no real-time interaction. | Algorithmic compensation: Employs corrective techniques (e.g., importance sampling) to mitigate latency effects. |
| Core Challenge | Exploration- exploitation trade-off during learning. *Resource utilization* under sequential tasks. | Distributional shift and limited data coverage. | Algorithmic stability under dynamic delays: High/variable latency causes importance weight variance explosion. |

F THE COMPARISON OF DIFFERENT REINFORCEMENT LEARNING PARADIGMS

F.1 Core Challenge: Off-Policy Learning and Importance Sampling

This formulation highlights the central challenge in asynchronous RL: the mismatch between the behavior policy (π_{θ_k}) and the target policy $(\pi_{\theta_{k+\tau}})$, which grows with τ and introduces bias and variance into the learning process. Addressing this mismatch under high and uncertain latency is the primary focus of our work.

The HeteroRL framework decouples rollout sampling from parameter learning, leading to a setting where the learner updates the new policy $\pi_{\theta_k+\tau}$ using data generated by an older policy π_{θ_k} —a canonical off-policy learning scenario. To correct for the distributional shift between the behavior and target policies, importance sampling (IS) is commonly employed. Under the GRPO framework, this weight is used to scale the advantage during policy gradient updates. However, when the delay τ is large, causing significant divergence between the policies (i.e., high KL divergence $D_{\text{KL}}(\pi_{\theta_{k+\tau}} || \pi_{\theta_k})$), we observe in practice that the variance of the importance weights increases rapidly with τ , and the estimation error of the expected reward also grows. As a result, the variance of w_{std} explodes, introducing high noise into gradient estimates and ultimately leading to training instability or collapse. The central challenge of this work is thus to design a robust algorithm that mitigates the training instability caused by the explosion of importance sampling variance under high network latency.

G CASE STUDY

A case of MATH level 3-5: GRPO v/s GSPO v/s GEPO

User:

Oser

Find the greatest integer less than $(\sqrt{7} + \sqrt{5})^6$. (Do not use a calculator!).

GRPO

<think> Okay, so I need to find the greatest integer less than $(\sqrt(7) + \sqrt(5))^6$. Hmm, that seems pretty complicated at first glance. Let me think about how to approach this. Wait, maybe ...(Omitted numerous ineffective thoughts and reflections)
Alternatively, maybe I can find a recurrence relation for S_n . Let me try to find a recurrence. Let me note that x and y are roots of the equation $t^2 - 2\sqrt{7}t + 2 = 0$. Therefore, for any $n \ge 2$, we have

$$x^n = 2\sqrt{7}x^{n-1} - 2x^{n-2}.$$

Similarly for y^n . Therefore,

$$S_n = x^n + y^n = 2\sqrt{7}(x^{n-1} + y^{n-1}) - 2(x^{n-2} + y^{n-2}) = 2\sqrt{7}S_{n-1} - 2S_{n-2}.$$

(This is the key step to solving this problem) ...

But this seems to involve $\sqrt{(7)}$, which complicates things. However, if I consider S_n + something, maybe I can eliminate the $\sqrt{(7)}$ term? Wait, but maybe I need another approach.

Alternatively, maybe I can compute S_n for n = 0, 1, 2, ... and find a pattern. Let me try that.

First, let's compute $S_0 = x^0 + y^0 = 1 + 1 = 2$.

(Numerous repetitions occur here.)

...(Omitted numerous ineffective thoughts and reflections)

Wait, but S_6 is $x^6 + y^6$. However, since $x = \sqrt(7) + \sqrt(5)$ and $y = \sqrt(7) - \sqrt(5)$, and since x > 0 and y is positive? Wait, $\sqrt(7)$ is about 2.645, $\sqrt(5)$ is about 2.236. So $\sqrt(7) + \sqrt(5)$ is about 4.881, and $\sqrt(7) - \sqrt(5)$ is about 0.409. So y is positive, but small. Therefore, x^6 is a large number, and y^6 is a small number. Therefore, $S_6 = x^6 + y^6$ is approximately x^6 , and the integer part of x^6 would be $S_6 - 1$, since y^6 is positive but less than 1. Therefore, the greatest integer less than x^6 is $S_6 - 1$.

But wait, let me check if y^6 is less than 1. Since $y = \sqrt(7) - \sqrt(5) \approx 0.409$, then y^6 is $(0.409)^6$. Let me compute that: $0.409^2 \approx 0.167$, then $0.167^3 \approx 0.004$

(Over length here, although close to arriving at the correct answer)

GSPO

<think> Okay, so I need to find the greatest integer less than $(\sqrt{7}) + \sqrt{5})^6$. Hmm, that seems a bit intimidating at first, but I'll try to break it down step by step. First, I remember that expressions like $(a+b)^n$ often have a pattern when expanded, especially when a and b are irrational numbers. Since the question is about the greatest integer less than this value, I might need to find an expression that's just below the actual value.
Let me think. If I consider $(\sqrt{7}) + \sqrt{5})^6$ and $(\sqrt{7}) - \sqrt{5})^6$, their sum might be an integer because the irrational parts cancel out. That's a common trick with conjugates. So, the sum of these two would be $(sqrt(7) + sqrt(5))^6 + (sqrt(7) - sqrt(5))^6$, which is an integer. Therefore, the integer part of the original expression would be that integer minus the $(sqrt(7) - sqrt(5))^6$, which is a small positive number.
So, the answer would be boxed as the boxed integer boxedboxedboxedboxed boxed boxed.

but boxed in the final answer boxed.

Therefore, the boxed answer is the boxed boxed

GEPO

<think> Okay, so I need to find the greatest integer less than $(\sqrt(7) + \sqrt(5))^6$. Hmm, that seems pretty big. I don't have a calculator, so I have to figure this out without calculating the actual value. Let me think.

••

But maybe I can find a recurrence formula. Let's see:

We know that x and y are roots of the quadratic equation $t^2 - (x + y)t + xy = 0$. So $t^2 - 2\sqrt{7}t + 2 = 0$. Therefore, the recurrence relation for S_n is $S_n = 2\sqrt{7}S_{n-1} - 2*S_{n-2}$. ($\mbox{$\frac{1}{2}$}$ This is the key step to solving this problem) ...

Therefore, the answer is 13535

1307 </think>

To find the greatest integer less than $(\sqrt{7}+\sqrt{5})^6$, we use the identity that for any expression of the form $(\sqrt{a}+\sqrt{b})^n+(\sqrt{a}-\sqrt{b})^n$, the result is an integer. This is because the irrational parts cancel out when the two expressions are added together.

Let: $-x = \sqrt{7} + \sqrt{5} - y = \sqrt{7} - \sqrt{5}$

We compute:

$$x^{6} + y^{6}$$

First, calculate x^2 :

$$x^{2} = (\sqrt{7} + \sqrt{5})^{2} = 7 + 2\sqrt{35} + 5 = 12 + 2\sqrt{35}$$

Then $x^4 = (x^2)^2 = (12 + 2\sqrt{35})^2 = 144 + 48\sqrt{35} + 4 \cdot 35 = 284 + 48\sqrt{35}$ Next, $x^6 = x^4 \cdot x^2 = (284 + 48\sqrt{35})(12 + 2\sqrt{35})$. Expanding this:

$$284 \cdot 12 + 284 \cdot 2\sqrt{35} + 48\sqrt{35} \cdot 12 + 48\sqrt{35} \cdot 2\sqrt{35}$$
$$= 3408 + 568\sqrt{35} + 576\sqrt{35} + 96 \cdot 35 = 3408 + 1144\sqrt{35} + 3360$$
$$= 6768 + 1144\sqrt{35}$$

Similarly, $y^6 = 6768 - 1144\sqrt{35}$ Adding x^6 and y^6 :

1323 Adding x^6 and y^6

$$x^6 + y^6 = 2 \cdot 6768 = 13536$$

Since $x^6 = 13536 - y^6$, and y^6 is a small positive number (as $\sqrt{7} - \sqrt{5} \approx 0.4097$), we have:

$$x^6 \approx 13536 - 0.004 \approx 13535.996$$

Thus, the greatest integer less than x^6 is:

Comparison:

GRPO exhibits a large number of ineffective reflections, wasting many valuable tokens; statistics show that GRPO uses words such as "wait" and "however" more than 22 times. Due to reward collapse, GSPO begins generating extensive repetitions after brief reasoning, rendering the response unreadable. GEPO arrives at the correct answer during the thinking phase through sound reasoning and reflection. When formally responding, its thought process is clear and steps are concise.

H FUTURE WORK

H.1 Defensive Sampling and Smooth Denominator Mechanism

As a promising direction for future work, we propose to explore a defensive sampling strategy that adaptively blends the target policy probability into the importance weight denominator. This approach aims to mitigate bias introduced by approximating the denominator while improving training stability—particularly in asynchronous or heterogeneous learning settings. By introducing a variance-aware smoothing coefficient, the method could dynamically interpolate between standard importance weighting and the policy gradient objective, thereby reducing sensitivity to high-variance estimates. We hypothesize that such a smooth,

adaptive weighting scheme would lead to more robust and stable policy updates, and plan to investigate its theoretical properties and empirical effectiveness in future studies.

I LARGE LANGUAGE MODEL USAGE

The use of LLMs in this article is limited to text polishing and code generation for plotting.