

# Exploiting Tree Structure for Credit Assignment in Reinforcement Learning with Large Language Models

Hieu Tran<sup>\*1,2</sup>, Zonghai Yao<sup>\*1,2</sup>, Hong Yu<sup>1,2,3</sup>

<sup>1</sup>Center for Healthcare Organization and Implementation Research, VA Bedford Health Care

<sup>2</sup>Manning College of Information and Computer Sciences, University of Massachusetts Amherst

<sup>3</sup>Miner School of Computer and Information Sciences, University of Massachusetts Lowell

## Abstract

Reinforcement learning has shown strong promise for strengthening the reasoning ability of large language models (LLMs), but sparse, delayed rewards over long chains make token-level credit assignment a central challenge. Actor-critic methods like PPO provide token-level credit but require training a value network alongside the policy, which introduces complexity and can encourage overfitting. Critic-free alternatives such as GRPO avoid this burden but rely on sequence-level outcomes, distributing a single reward uniformly across tokens and ignoring structural differences between responses. We propose Prefix-to-Tree (P2T), which organizes the sampled responses of a prompt into a prefix tree and computes nonparametric prefix values by aggregating descendant outcomes. Building on this idea, we develop TEMPO (Tree-Estimated Mean Prefix Value for Policy Optimization), a critic-free algorithm that enriches GRPO with branch-aware temporal-difference (TD) corrections. Across Qwen3-1.7B and Qwen3-4B, TEMPO consistently improves both convergence and final performance over PPO and GRPO on in-distribution benchmarks (MATH, MedQA) and out-of-distribution settings (GSM-HARD, AMC23, MedMCQA, MMLU-Medical), achieving higher validation accuracy within comparable wall-clock time. <sup>1</sup>

## 1 Introduction

Reinforcement learning (RL) (Sutton et al., 1998) is an effective way to strengthen the reasoning of large language models (LLMs) (Zhang et al., 2025). In LLM settings, rewards are sparse and delayed and sequences are long (Jaech et al., 2024; Guo et al., 2025), so the key challenge is **credit assignment**: give the outcome reward to the few tokens that really change the solution. We study the

*verifiable-reward* setting, where the final answer for a prompt is checkable and we can draw multiple responses for the same prompt. This is common in long “thinking” or chain-of-thought (CoT) tasks such as mathematics and medical QA, where most steps are low-impact and only a small set of *decision tokens* (e.g., strategy choice, formula selection, diagnostic commitment) moves the outcome. Aggregating multiple responses naturally induces an *implicit prefix tree*: internal nodes are shared prefixes and branch nodes mark decision points with multiple plausible continuations. A good learning rule should use the branching structure across responses and focus credit on those decision points.

**Proximal Policy Optimization (PPO)** gives token-level advantages with a learned value and generalized advantage estimation (GAE), which mixes Monte Carlo (MC) returns with temporal-difference (TD) bootstrapping (Schulman et al., 2015, 2017). However, jointly training the actor and critic is complex and often fails to generalize, as critic-derived token-level values can induce overfitting (Wang et al., 2025b; Chaudhari et al., 2024). **Group Relative Policy Optimization (GRPO)** removes the critic and uses group-relative baselines over responses to the same prompt (Shao et al., 2024; Yu et al., 2025). It is simple and fits verifiable rewards. However, it spreads a single sequence-level signal across all tokens and overlooks mid-trajectory decisions. As a result, token-level credit is weak when reasoning branches. Recent “key-token ideas” (Wang et al., 2025a) move toward finer signals by focusing gradient updates on high-entropy tokens. This approach benefits from exploiting response structure and concentrating learning on decision-heavy positions. However, while entropy-based updates can improve exploitation within known reasoning patterns, they may struggle to acquire new domain knowledge where exploration across broader patterns of the responses is necessary.

Equal contribution

<sup>1</sup>Our code can be accessed at: <https://github.com/fatebreaker/tempo>

We present **Prefix-to-Tree (P2T)** as a simple procedure that converts a group of responses into a prefix tree and computes *nonparametric* prefix values  $V(s)$  by averaging descendant returns. Building on P2T, we introduce **TEMPO (Tree-Estimated Mean Prefix Value for Policy Optimization)**, a critic-free policy optimization method that restores token-level credit only where it matters. For each prompt, sampled responses form paths in the implicit prefix tree. TEMPO augments the group-relative outcome signal of GRPO with *branch-gated* temporal-difference (TD) corrections derived from the tree: at non-branching tokens  $V(s_{t+1}) = V(s_t)$  and the TD error is zero, so the update reduces to GRPO, while at branching tokens it supplies precise token-level credit. TEMPO maintains the GRPO training loop and cost. It does not train a value model or add a process reward model or a judge, it also does not require a teacher or a new sampler.

**Empirical scope and applicability.** Across *Qwen3-1.7B* and *Qwen3-4B*, TEMPO consistently attains higher accuracy than PPO (Schulman et al., 2017), GRPO (Shao et al., 2024), HEPO (Wang et al., 2025a), and TREERPO (Yang et al., 2025b) on both in-distribution (MATH, MedQA) and out-of-distribution (GSM-HARD, AMC23, MedMCQA, MMLU-Medical) evaluations, while reaching strong validation performance in less wall-clock time. Validation curves indicate that, on math reasoning, approaches that emphasize token-level structure, such as HEPO (Wang et al., 2025a) (e.g., focusing updates on high-entropy decision tokens) already enjoy an advantage, suggesting the RL phase mainly reinforces reasoning patterns learned during pretraining and SFT. Yet, TEMPO goes further by injecting tree-gated TD credit at the exact branching points. On medical reasoning, where domain knowledge must be newly acquired, methods that rely on group-relative exploration such as GRPO (Shao et al., 2024) generalize better than purely exploitation-oriented updates; TEMPO combines this robust group baseline with branch-aware TD from P2T’s nonparametric prefix values, improving both convergence speed and final accuracy. In practice, TEMPO is most beneficial when rewards are verifiable and prompts yield meaningful branching, delivering precise token-level credit without a value network or auxiliary judges, and serving as a drop-in, efficiency-preserving upgrade to GRPO-style training. Overall, our key **contributions** include:

1. We introduce **Prefix-to-Tree (P2T)**, a simple procedure that converts each prompt’s group of responses into a prefix tree and derives *non-parametric* prefix values  $V(s_t)$  by aggregating descendant outcomes.
2. Building on P2T, we propose **TEMPO**, a drop-in, GRPO-compatible algorithm that augments the group-normalized outcome signal with *branch-gated* TD corrections, providing precise token-level credit at decision points while retaining GRPO-like compute and simplicity.
3. On *Qwen3-1.7B/4B*, TEMPO improves convergence speed and final accuracy over other baselines on in-distribution (MATH, MedQA) and out-of-distribution (GSM-HARD, AMC23, MedMCQA, MMLU-Medical) benchmarks under the same hardware budget.

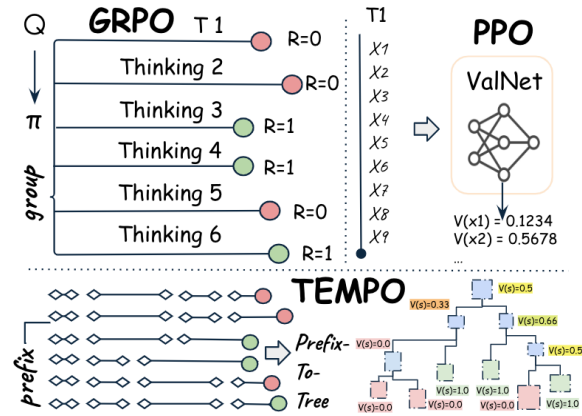


Figure 1: **Comparison of credit assignment for RL training with verifiable rewards.** GRPO: all tokens in each sampled answer share one sequence-level return; branching is ignored so credit spreads evenly. PPO: a learned value network estimates  $V(s_t)$  and provides token-level advantages via GAE, but requires a critic and higher compute. TEMPO: convert the answer group for one prompt into a prefix tree and compute *nonparametric* prefix values  $V(s)$  by averaging descendant outcomes; use *branch-gated* TD corrections to assign credit at branches.

## 2 Related Work

Credit assignment is central in post-training for reasoning LLMs. RLHF brought PPO with a learned value (critic) and GAE to reduce variance (Ouyang et al., 2022; Schulman et al., 2017, 2015). This improved alignment, however, comes at the cost of

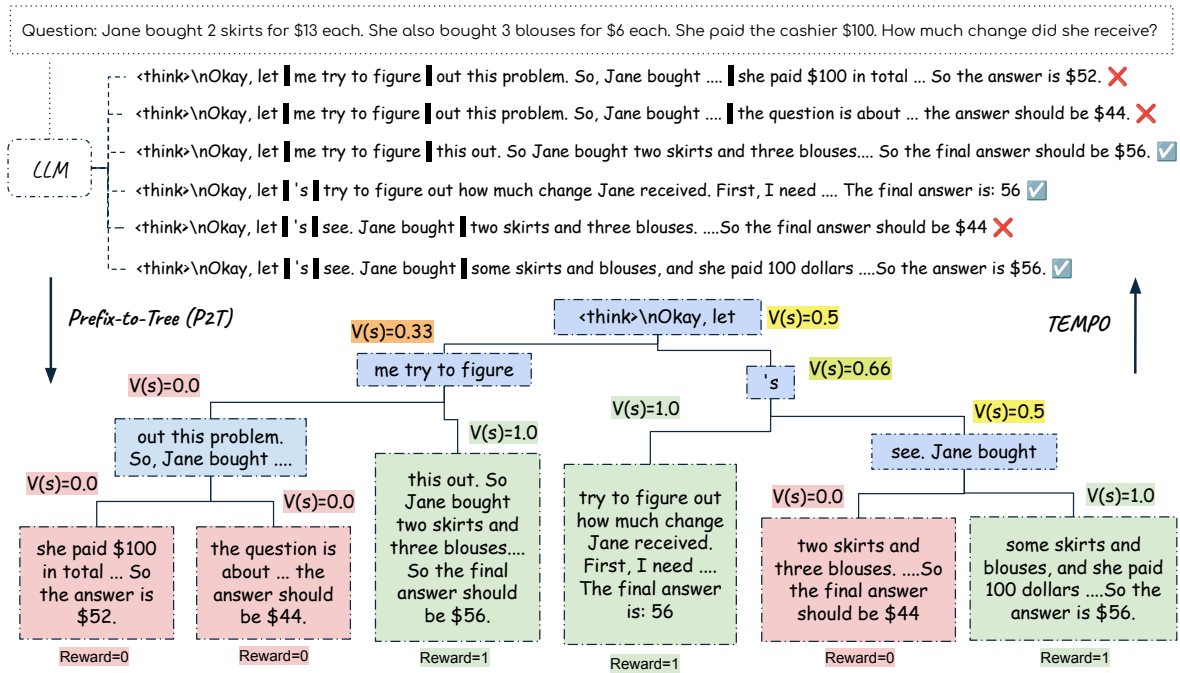


Figure 2: Overview of prefix tree value estimation in TEMPO. Each node corresponds to a token prefix  $s$ , with  $V(s)$  estimated by averaging over the outcomes of all descendant completions. Green leaves denote correct responses ( $r = 1$ ), red leaves denote incorrect ones ( $r = 0$ ). Intermediate nodes inherit averaged values (e.g.,  $V(s) = 0.5$ ), providing informative signals at branching points.

critic training, which adds complexity and tuning, and value prediction is brittle on long chains. To avoid a critic, several lines move towards value-free or RL-free updates that treat the entire response as a single action. DPO optimizes pairwise preferences in an offline bandit view (Rafailov et al., 2023). Rejection-sampling methods, such as RestEM, fine-tune only on full high-reward responses (Singh et al., 2023). RLOO, GRPO, and DAPO compute group-normalized sequence advantages over multiple samples of the same prompt, thereby removing the value network (Ahmadian et al., 2024; Shao et al., 2024; Yu et al., 2025). These methods are simple and stable, but their feedback is sequence-level and credits all tokens equally, which weakens token-level credit in long reasoning. VinePPO instead replaces the critic with Monte Carlo estimates obtained by re-sampling continuations from each text prefix, yielding accurate prefix values in language settings (Kazemnejad et al., 2024). However, it requires fresh rollouts at many branch nodes and raising sampling cost when trees are wide or branch early and still uses path-wise PPO advantages without group-normalized baselines.

A second thread tries to push feedback below the sequence. Token or span-level preference and dense-reward methods give finer signals (Yoon

et al., 2024; Yang et al., 2024; Chan et al., 2024). Process supervision verifies intermediate steps or chain consistency to localize the first error (Lightman et al., 2023; Chen et al., 2024b; Setlur et al., 2024b; Chen et al., 2024a; Zhang et al., 2024b). Yet many step-by-step or tree-style approaches depend on a learned process reward model (PRM) or a judge to score nodes, which re-introduces a value function and adds verifier training cost. Some recent work also injects ad-hoc Monte Carlo (MC) signals into DPO to flag faulty steps (Hwang et al., 2024; Setlur et al., 2024a). Our approach follows the value-free direction but uses the trajectory’s structure: it forms non-parametric prefix values from sibling continuations within a prompt group, then applies temporal-difference updates only at branching tokens where returns diverge, while non-branch tokens fall back to a GRPO-like baseline. In this way the update is simple like GRPO, yet it concentrates credit on decision points without a PRM.

Several contemporaneous works exploit the *tree structure* of rollouts to densify credit assignment and/or cut sampling cost. TreePO (Li et al., 2025) reframes on-policy rollouts as a tree search with segmented decoding and heuristic branching/fall-back, amortizing shared prefixes (KV caching) and

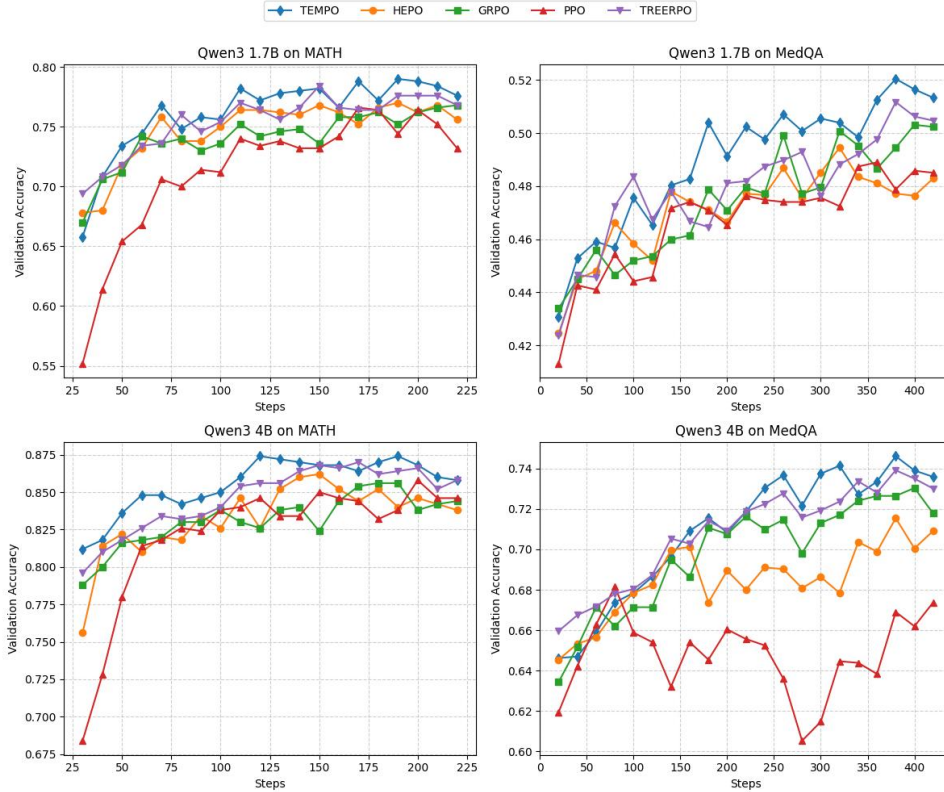


Figure 3: Validation accuracy of MATH and MedQA for Qwen3-1.7B and Qwen3-4B. We compare TEMPO with PPO, GRPO, and HEPO. TEMPO consistently achieves higher accuracy and faster convergence across both domains and model sizes.

introducing a tree-based segment-level advantage estimator; this improves stability and reduces sampling compute while maintaining or improving accuracy. TREERPO (Yang et al., 2025b) extends GRPO by performing explicit tree sampling and forming step-level sibling groups to estimate expected rewards per step, yielding dense process signals and reporting consistent gains over GRPO with shorter responses. TreeRL (Hou et al., 2025) integrates an entropy-guided sampler (EPTree) that branches at uncertain tokens, then back-propagates leaf rewards to provide global and local (step) advantages thereby eliminating a separate process reward model. Tree-OPO (Huang et al., 2025) leverages *off-policy* teacher MCTS to build prefix trees and proposes staged, prefix-conditioned advantage estimation to stabilize GRPO-style updates. Unlike methods that require dedicated tree samplers (TreeRPO/TreeRL) or off-policy teacher trees (Tree-OPO), TEMPO operates in the standard GRPO setting and treats the *implicit* prefix tree formed by a group of responses as a nonparametric value baseline: it computes  $V(s_t)$  from all completions sharing the prefix  $s_t$  and adds a token-

level TD correction to the group-relative (Monte Carlo) signal. This yields branch-aware advantages without a learned value network, extra reward/process models, or special search procedures, while remaining fully on-policy and drop-in compatible with GRPO training loops.

### 3 Preliminaries

We begin by reviewing the advantage estimation used in Proximal Policy Optimization and Group Relative Policy Optimization.

**PPO.** PPO (Schulman et al., 2017) is a policy gradient method that stabilizes updates via a clipped objective. A key component is the estimation of the advantage function  $A_t$ , which measures how much better an action  $a_t$  is compared to the average action at state  $s_t$ . PPO commonly employs *generalized advantage estimation* (GAE) (Schulman et al., 2017), defined as

$$\hat{A}_t^{\text{GAE}(\gamma, \lambda)} = \sum_{l=0}^{T-t-1} (\gamma\lambda)^l \delta_{t+l}$$

In the original formulation,  $\gamma$  serves as a dis-

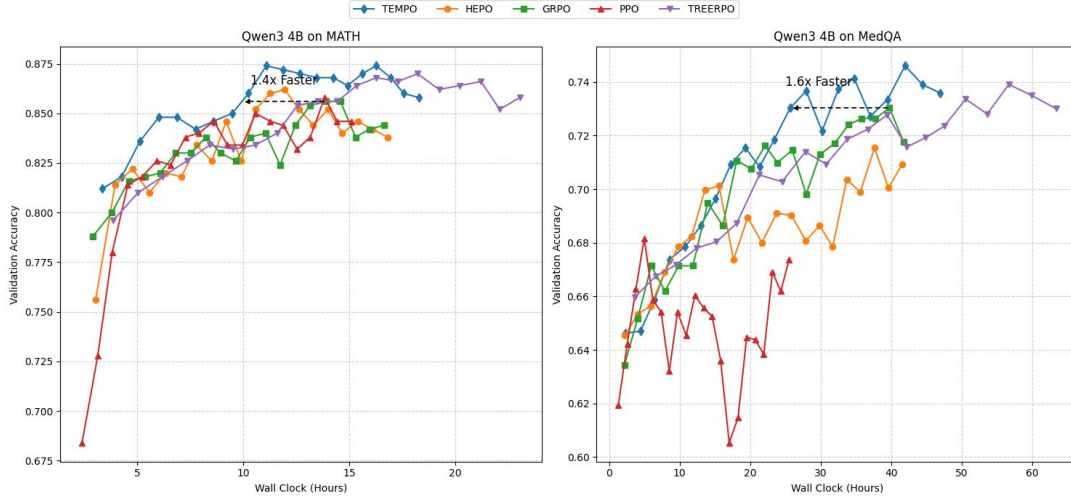


Figure 4: TEMPO converges faster and to higher accuracy than GRPO, passes GRPO’s peak performance in fewer iterations and less overall time.

count factor that reduces the weight of delayed rewards and helps stabilize infinite-horizon settings. However, in the context of large language model (LLM) training, it is common to set  $\gamma = 1.0$  so that long completions are not penalized relative to short ones. With this setting, the GAE formula simplifies to

$$\hat{A}_t^{\text{GAE}(\lambda)} = \sum_{l=0}^{T-t-1} \lambda^l \delta_{t+l}$$

The bias–variance tradeoff is then controlled solely by the parameter  $\lambda$ .

**Special cases.**

- $\lambda = 0$  (*TD(0)*).

$$\hat{A}_t^{\lambda=0} = \delta_t = r_t + V(s_{t+1}) - V(s_t),$$

the one-step temporal-difference error (lowest variance, highest bias).

- $\lambda = 1$  (*Monte Carlo*).

$$\hat{A}_t^{\lambda=1} = \sum_{l=0}^{T-t-1} r_{t+l} + V(s_T) - V(s_t).$$

If  $s_T$  is terminal so  $V(s_T) = 0$ , then

$$\hat{A}_t^{\lambda=1} = \left( \sum_{l=0}^{T-t-1} r_{t+l} \right) - V(s_t),$$

i.e., the full Monte Carlo return minus the baseline (unbiased, higher variance).

**GRPO.** GRPO (Shao et al., 2024) was designed for reinforcement learning with verifiable feedback. Formally, for each question  $q$ , a group of  $G$  responses  $\{o_1, \dots, o_G\}$  is sampled from the old policy  $\pi_{\theta_{\text{old}}}$ , and a reward model assigns scores  $r = \{r_1, \dots, r_G\}$ . These rewards are then normalized by subtracting the group mean and dividing by the group standard deviation. Under outcome supervision, the normalized reward  $\tilde{r}_i = \frac{r_i - \text{mean}(r)}{\text{std}(r)}$  is applied uniformly to all tokens of output  $o_i$ , so that

$$\hat{A}_{i,t} = \tilde{r}_i = \frac{r_i - \text{mean}(r)}{\text{std}(r)}, \quad \forall t \in o_i.$$

Under process supervision, the same normalization is applied at the step level, and the normalized step rewards are distributed to the corresponding tokens. GRPO thus avoids training a separate value model and provides efficient group-relative baselines, but it relies purely on Monte Carlo outcomes and discards token-level temporal structure. This setup highlights the gap: PPO provides token-level advantages but requires a value model, while GRPO is model-free but trajectory-level only. Our method, TEMPO, combines the strengths of both.

## 4 Methodology

### 4.1 Value Estimation from Prefix Tree

Figure 2 illustrates how TEMPO derives value estimates directly from the tree structure formed by a group of sampled responses. Each path in the tree corresponds to a response generated by the policy, and each node represents a token prefix  $s_t$

up to time  $t$ . The tree branches whenever different responses diverge at a given token. Terminal nodes are assigned rewards  $r \in \{0, 1\}$  based on verifiable correctness (e.g., whether the final answer matches the ground truth).

Instead of training a separate value model as in PPO, TEMPO computes  $V(s_t)$  directly from the group of trajectories. For a given prefix  $s_t$ , the value is estimated as the average normalized reward of all descendant completions that share this prefix:

$$V(s_t) = \frac{1}{|D(s_t)|} \sum_{j \in D(s_t)} r_j,$$

where  $D(s_t)$  is the set of responses passing through  $s_t$ , and  $r_j$  is the outcome reward. This provides a *value function* without introducing an additional learned critic.

In the example shown in Figure 2, some prefixes lead to correct answers ( $r = 1$ ) while others lead to incorrect ones ( $r = 0$ ). TEMPO propagates these signals upward by averaging over the subtree, yielding intermediate values (e.g.,  $V(s_t) = 0.5$  when half of the descendant completions are correct). As a result, branch nodes obtain informative value estimates that reflect the quality of their continuations, while non-branch nodes naturally inherit consistent values from their unique continuation. This design ensures that tokens along successful reasoning paths (green leaves) contribute positively to the estimated value of their prefixes, tokens along failed reasoning paths (red leaves) reduce the value of their prefixes and branching points receive *discriminative signals*, as the value function captures how sibling continuations differ in correctness. By computing  $V(s_t)$  directly from the tree, TEMPO provides token-level evaluative feedback while remaining model-free, combining the efficiency of GRPO with the structured credit assignment of PPO.

## 4.2 Branch-Aware Advantage Estimation

Having defined the prefix tree value function  $V(s_t)$ , we now describe how TEMPO constructs advantages by combining *response-level Monte Carlo signals* from GRPO with *token-level temporal-difference corrections* derived from the tree.

**Response-level (MC) signal.** GRPO provides outcome-level supervision by normalizing the rewards across a group of  $G$  responses. For outcome supervision, each response  $o_i$  receives a normalized

reward

$$\tilde{r}_i = \frac{r_i - \text{mean}(r)}{\text{std}(r)},$$

and assigns it uniformly to all tokens of the trajectory. This yields a pure Monte Carlo signal: every token in a response inherits the same scalar advantage  $\tilde{r}_i$ . While efficient, this discards the structure of reasoning trajectories.

**Token-level (TD) correction.** TEMPO augments this outcome-level signal with a token-level TD term based on branch-aware values. For token  $t$  in trajectory  $i$ , with state prefix  $s_t$  and successor  $s_{t+1}$ , we define the TD error as

$$\delta_{i,t} = V(s_{t+1}) - V(s_t).$$

This term captures how much the estimated value changes when extending from prefix  $s_t$  to  $s_{t+1}$ . Importantly,  $\delta_{i,t}$  is only nonzero at branching points, since non-branch tokens have identical descendant outcomes and thus  $V(s_{t+1}) = V(s_t)$ .

**Combined TEMPO advantage.** The final TEMPO advantage integrates both levels of signal:

$$\hat{A}_{i,t} = \frac{1}{\text{std}(r)} \underbrace{[r_i - \text{mean}(r)]}_{\text{MC signal}} + \underbrace{[V(s_{t+1}) - V(s_t)]}_{\text{TD error}}$$

The *MC signal* provides global outcome-level supervision, aligning each response relative to its group. The *TD error* provides local, branch-aware token-level feedback, highlighting where reasoning paths diverge in quality.

## Interpretation.

- At non-branch tokens,  $\delta_{i,t} = 0$  and TEMPO reduces to GRPO.
- At branch tokens,  $\delta_{i,t}$  propagates differences in descendant correctness back to the branching point, allowing finer credit assignment.
- Unlike PPO, TEMPO does not train a separate value model:  $V(s_t)$  is estimated directly from the group of responses.

## 4.3 Policy Update

For the policy update, TEMPO follows the same principles as DAPO (Yu et al., 2025), incorporating several practical design choices that improve stability and efficiency such as Clip-Higher (decoupled clipping), token-level policy-gradient loss (global token averaging) and remove KL divergence.

**TEMPO loss function.** Combining these practices with our proposed branch-aware advantage estimation, the loss is defined as

$$\mathcal{J}_{\text{TEMPO}}(\theta) = \mathbb{E}_{q, \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|q)} \left[ \frac{1}{\sum_i |o_i|} \sum_{i=1}^G \sum_{t=1}^{|o_i|} \min \left( r_{i,t}(\theta) \hat{A}_{i,t}, \text{clip}(r_{i,t}(\theta), 1 - \varepsilon_{\text{low}}, 1 + \varepsilon_{\text{high}}) \hat{A}_{i,t} \right) \right] \quad (1)$$

## 5 Experimental Setup

**Datasets and Models** We train on MATH (Hendrycks et al.) and MedQA (Jin et al., 2021) as in-distribution tasks, and evaluate on their test sets plus OOD benchmarks (GSM-HARD, AMC23, MedMCQA, MMLU-Medical). Experiments use Qwen3-1.7B and Qwen3-4B (Yang et al., 2025a) under identical settings (details in Appendix A.2).

**Baselines** Our main baseline is GRPO (Shao et al., 2024), which incorporates several practical strategies from DAPO (Yu et al., 2025): removing the KL penalty, introducing a clip-higher mechanism, and applying a token-level policy gradient loss. These modifications make GRPO one of the state-of-the-art RLVF algorithms without requiring a value network. We further include TREERPO (Yang et al., 2025b), which constructs sampled trees and computes step-level, group-relative rewards to provide denser credit signals without a separate reward model (Yang et al., 2025b). TEMPO builds on GRPO and improves credit assignment by exploiting the tree structure of responses. We also compare against a GRPO variant that targets *high-entropy minority tokens* (Wang et al., 2025a), where gradient updates are applied only to high-entropy tokens. For clarity in experiments and figures, we denote this variant as HEPO (High Entropy Policy Optimization). Finally, we include an actor-critic baseline: PPO (Schulman et al., 2017), where the critic model is matched in size to the actor model.

## 6 Results

In this section, we evaluate the effect of better Credit Assignment on task performance, efficiency, and generalization dynamics.

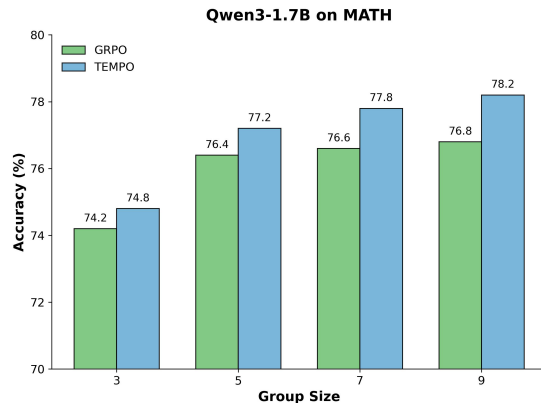


Figure 5: Effect of group size on MATH accuracy for Qwen3-1.7B. TEMPO consistently outperforms GRPO across all settings.

### 6.1 Task Performance

Figure 3 shows validation curves on MATH and MedQA for both Qwen3-1.7B&4B. Across settings, TEMPO has the best performance in terms of convergence speed and final accuracy. On MATH, TEMPO consistently outperforms others, followed by TREERPO, then HEPO, GRPO, and PPO. The fact that HEPO performs slightly better than GRPO and PPO suggests that focusing updates on high-entropy tokens helps exploit the reasoning structures already present in the model. Since mathematical reasoning knowledge is largely captured during pretraining and instruction tuning, the RL process primarily reinforces existing structures rather than learning new ones. Thus, token-structure-oriented methods, such as HEPO, gain an advantage compared with GRPO. On MedQA, however, the trend differs. TEMPO again delivers the best results, but GRPO surpasses HEPO, and PPO lags behind. We hypothesize that medical reasoning requires learning novel domain-specific knowledge, which emphasizes *exploration* rather than pure exploitation of existing token-level structures. Here, GRPO’s group-relative normalization provides stronger signals than PPO, while HEPO’s focus on high-entropy tokens is insufficient to capture new knowledge. TEMPO combines the benefits of GRPO with tree-structured TD guidance, enabling effective exploration while still leveraging structural signals, leading to the best generalization in the medical domain.

### 6.2 Computational Efficiency

Figure 4 reports validation accuracy against wall-clock training time on Qwen3-4B for both MATH



Model	Math			Medical		
	MATH	GSM-HARD	AMC23	MedQA	MedMCQA	MMLU-Medical
Qwen3-1.7B	68.5	46.85	57.5	46.11	43.17	57.85
+ PPO	81.6	53.37	67.5	52.94	48.05	70.16
+ GRPO	82.4	53.15	72.5	56.24	49.56	71.53
+ HEPO	81.7	52.09	62.5	54.28	48.98	71.35
+ TREERPO	84.6	54.82	72.5	57.42	50.23	72.08
+ TEMPO	<b>87.0</b>	<b>57.69</b>	<b>75.0</b>	<b>59.15</b>	<b>51.54</b>	<b>73.37</b>
Qwen3-4B	71.3	54.13	75.0	65.36	56.63	78.24
+ PPO	87.4	58.07	85.0	72.03	59.29	83.10
+ GRPO	87.6	59.81	85.0	76.12	60.55	83.19
+ HEPO	88.2	59.51	82.5	74.31	59.48	82.37
+ TREERPO	88.6	60.27	87.5	77.06	60.96	83.56
+ TEMPO	<b>91.0</b>	<b>62.32</b>	<b>92.5</b>	<b>79.18</b>	<b>62.51</b>	<b>85.49</b>

Table 1: Comparison of PPO, GRPO, HEPO, and TEMPO on mathematical and medical reasoning benchmarks using Qwen3-1.7B and Qwen3-4B as base models. MATH and MedQA are considered *in-distribution* (ID) tasks, while GSM-HARD, AMC23, MedMCQA, and MMLU-Medical are treated as *out-of-distribution* (OOD) evaluations.

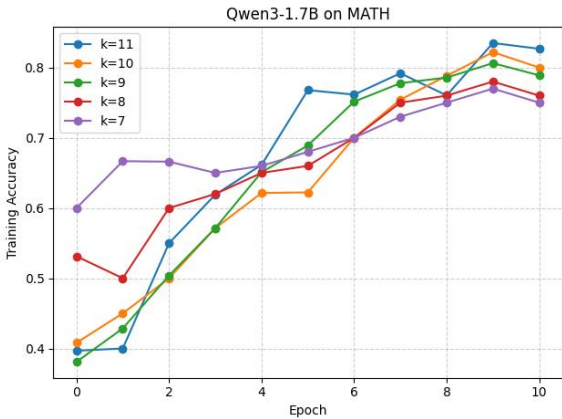


Figure 7: Training accuracy vs. epochs for different numbers of preserved branches  $k$  at group size  $G=7$ .

benchmarks for mathematics and medicine. On the ID tasks (MATH and MedQA), TEMPO consistently achieves the highest accuracy, surpassing TREERPO, PPO, GRPO, and HEPO across both Qwen3-1.7B and Qwen3-4B. For example, TEMPO improves MedQA accuracy from 76.1% (GRPO, 4B) to 79.2%, and raises MATH accuracy from 87.6% to 91.0%. On OOD evaluations, TEMPO also establishes clear gains. In the math domain, it pushes GSM-HARD accuracy from 59.8% (GRPO, 4B) to 62.3%, and AMC23 from 85.0% to 92.5%. In the medical domain, it improves MedMCQA from 60.55% to 62.51% and MMLU-Medical from 83.2% to 85.5%. These improvements across unseen distributions high-

light that TEMPO not only enhances in-distribution learning efficiency but also yields stronger generalization to harder and more diverse reasoning tasks.

## 7 Conclusion

In this work, we introduced TEMPO, a reinforcement learning algorithm for LLM alignment that integrates TD signals into group-relative optimization by exploiting the tree structure of sampled responses. Unlike PPO, which requires training a separate value network, and GRPO, which discards token-level information by relying purely on Monte Carlo signals, TEMPO unifies the strengths of both approaches without additional model components. By deriving value estimates directly from the responses tree, TEMPO enables token-level TD corrections on top of group-relative normalization, yielding more fine-grained and stable credit assignment. Our experiments across mathematics and medicine demonstrate two key findings. First, TEMPO achieves higher accuracy than PPO, GRPO, HEPO and TREERPO on both in-distribution and out-of-distribution benchmarks, showing strong generalization. Second, TEMPO converges significantly faster in wall-clock time, achieving comparable or better accuracy under the same hardware configuration. Overall, TEMPO establishes a practical and scalable approach to reinforcement learning with verifiable feedback.

## 8 Limitations

While TEMPO demonstrates consistent improvements over PPO, GRPO, HEPO, and TREERPO across both mathematical and medical reasoning benchmarks, several limitations remain.

First, TEMPO assumes a verifiable-reward setting where final outcomes are checkable and multiple responses per prompt can be sampled. This restricts applicability in domains without objective verification signals or where response sampling is expensive.

Second, the effectiveness of branch-gated TD corrections depends on the existence of meaningful branching structure; when trajectories diverge early or remain largely linear, the additional signal diminishes and TEMPO behaves similarly to GRPO.

Third, our experiments focus on Qwen3-1.7B and Qwen3-4B and primarily on math and medical QA; future work is needed to assess scalability to larger models, broader domains, and multimodal reasoning.

Finally, while results suggest better generalization, we do not yet study robustness under adversarial inputs, noisy supervision, or deployment constraints. Exploring adaptive branching strategies, richer token-level credit signals, and extensions beyond verifiable rewards remain promising directions.

## 9 Ethics Statement

All authors have read and will adhere to the ICLR Code of Ethics. We acknowledge that the Code applies to all aspects of conference participation (submission, reviewing, discussion).

**Scope and compliance.** This work studies RL for LLM reasoning using public benchmarks only. No human subjects were recruited and no new personal data were collected; therefore, IRB approval was not required.

**Data, privacy, and licensing.** We use publicly available datasets (MATH, MedQA, GSM-HARD, AMC23, MedMCQA, MMLU-Medical) under their licenses. To the best of our knowledge, these datasets do not contain personally identifiable information. Evaluation relies on programmatic, verifiable feedback (e.g., exact match), not human raters or proprietary judges.

**Safety and misuse.** Models are research artifacts and *not* clinical decision tools. Medical

benchmarks are knowledge tests; deploying models downstream for healthcare or other high-stakes uses requires additional validation, domain oversight, and regulatory compliance.

**Fairness and limitations.** Benchmarks may encode societal or domain biases. We report results across multiple tasks and discuss limitations (e.g., distribution shift, overfitting risk). We encourage cautious interpretation beyond the evaluated settings.

**Transparency and reproducibility.** We document training settings, hardware, and evaluation protocols; we intend to release code/configs subject to third-party licenses. No hidden reward/process models, teachers, or special samplers were used. For details on LLM usage in paper writing, see Appendix A.8.

## References

- Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. 2024. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv preprint arXiv:2402.14740*.
- Alex J Chan, Hao Sun, Samuel Holt, and Mihaela Van Der Schaar. 2024. Dense reward for free in reinforcement learning from human feedback. *arXiv preprint arXiv:2402.00782*.
- Shreyas Chaudhari, Pranjal Aggarwal, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, Karthik Narasimhan, Ameet Deshpande, and Bruno Castro da Silva. 2024. Rlhf deciphered: A critical analysis of reinforcement learning from human feedback for llms. *ACM Computing Surveys*.
- Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai Fan. 2024a. Alphamath almost zero: process supervision without process. *Advances in Neural Information Processing Systems*, 37:27689–27724.
- Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai Fan. 2024b. Step-level value preference optimization for mathematical reasoning. *arXiv preprint arXiv:2406.10858*.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Pal: Program-aided language models. In *International Conference on Machine Learning*, pages 10764–10799. PMLR.
- Evan Greensmith, Peter L Bartlett, and Jonathan Baxter. 2004. Variance reduction techniques for gradient estimates in reinforcement learning. *Journal of Machine Learning Research*, 5(Nov):1471–1530.

- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- 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. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Zhenyu Hou, Ziniu Hu, Yujiang Li, Rui Lu, Jie Tang, and Yuxiao Dong. 2025. Treerl: Llm reinforcement learning with on-policy tree search. *arXiv preprint arXiv:2506.11902*.
- Bingning Huang, Tu Nguyen, and Matthieu Zimmer. 2025. Tree-ppo: Off-policy monte carlo tree-guided advantage optimization for multistep reasoning. *arXiv preprint arXiv:2509.09284*.
- Shengyi Huang, Michael Noukhovitch, Arian Hosseini, Kashif Rasul, Weixun Wang, and Lewis Tunstall. 2024. The n+ implementation details of rlhf with ppo: A case study on tl; dr summarization. *arXiv preprint arXiv:2403.17031*.
- Hyeonbin Hwang, Doyoung Kim, Seungone Kim, Seonghyeon Ye, and Minjoon Seo. 2024. Self-explore: Enhancing mathematical reasoning in language models with fine-grained rewards. *arXiv preprint arXiv:2404.10346*.
- Hamish Ivison, Yizhong Wang, Jiacheng Liu, Zeqiu Wu, Valentina Pyatkin, Nathan Lambert, Noah A Smith, Yejin Choi, and Hanna Hajishirzi. 2024. Unpacking dpo and ppo: Disentangling best practices for learning from preference feedback. *Advances in neural information processing systems*, 37:36602–36633.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. 2024. Openai o1 system card. *arXiv preprint arXiv:2412.16720*.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *Applied Sciences*, 11(14):6421.
- Amirhossein Kazemnejad, Milad Aghajohari, Eva Portelance, Alessandro Sordani, Siva Reddy, Aaron Courville, and Nicolas Le Roux. 2024. Vineppo: Refining credit assignment in rl training of llms. *arXiv preprint arXiv:2410.01679*.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th symposium on operating systems principles*, pages 611–626.
- Yizhi Li, Qingshui Gu, Zhoufutu Wen, Ziniu Li, Tianshun Xing, Shuyue Guo, Tianyu Zheng, Xin Zhou, Xingwei Qu, Wangchunshu Zhou, et al. 2025. Treepo: Bridging the gap of policy optimization and efficacy and inference efficiency with heuristic tree-based modeling. *arXiv preprint arXiv:2508.17445*.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Ankit Pal, Logesh Kumar Umapathi, and Malaikandan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering. In *Conference on health, inference, and learning*, pages 248–260. PMLR.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. 2015. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Amrith Setlur, Saurabh Garg, Xinyang Geng, Naman Garg, Virginia Smith, and Aviral Kumar. 2024a. RL on incorrect synthetic data scales the efficiency of llm math reasoning by eight-fold. *Advances in Neural Information Processing Systems*, 37:43000–43031.
- Amrith Setlur, Chirag Nagpal, Adam Fisch, Xinyang Geng, Jacob Eisenstein, Rishabh Agarwal, Alekh Agarwal, Jonathan Berant, and Aviral Kumar. 2024b. Rewarding progress: Scaling automated process verifiers for llm reasoning. *arXiv preprint arXiv:2410.08146*.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Avi Singh, John D Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, Xavier Garcia, Peter J Liu, James Harrison, Jaehoon Lee, Kelvin Xu, et al. 2023. Beyond human data: Scaling self-training for problem-solving with language models. *arXiv preprint arXiv:2312.06585*.

- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. 2023. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180.
- Richard S Sutton, Andrew G Barto, et al. 1998. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge.
- Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen, Jianxin Yang, Zhenru Zhang, et al. 2025a. Beyond the 80/20 rule: High-entropy minority tokens drive effective reinforcement learning for llm reasoning. *arXiv preprint arXiv:2506.01939*.
- Zihan Wang, Kangrui Wang, Qineng Wang, Pingyue Zhang, Linjie Li, Zhengyuan Yang, Xing Jin, Kefan Yu, Minh Nhat Nguyen, Licheng Liu, et al. 2025b. Ragen: Understanding self-evolution in llm agents via multi-turn reinforcement learning. *arXiv preprint arXiv:2504.20073*.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025a. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Kailai Yang, Zhiwei Liu, Qianqian Xie, Jimin Huang, Erxue Min, and Sophia Ananiadou. 2024. Selective preference optimization via token-level reward function estimation. *arXiv preprint arXiv:2408.13518*.
- Zhicheng Yang, Zhijiang Guo, Yinya Huang, Xiaodan Liang, Yiwei Wang, and Jing Tang. 2025b. Treerpo: Tree relative policy optimization. *arXiv preprint arXiv:2506.05183*.
- Eunseop Yoon, Hee Suk Yoon, Soohwan Eom, Gunsoo Han, Daniel Wontae Nam, Daejin Jo, Kyoungwoon On, Mark A Hasegawa-Johnson, Sungwoong Kim, and Chang D Yoo. 2024. Tlcr: Token-level continuous reward for fine-grained reinforcement learning from human feedback. *arXiv preprint arXiv:2407.16574*.
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. 2025. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*.
- Chi Zhang, Guangming Sheng, Siyao Liu, Jiahao Li, Ziyuan Feng, Zherui Liu, Xin Liu, Xiaoying Jia, Yanghua Peng, Haibin Lin, et al. 2024a. A framework for training large language models for code generation via proximal policy optimization. In *NL2Code Workshop of ACM KDD*.
- Dan Zhang, Sining Zhoubian, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. 2024b. Rest-mcts\*: Llm self-training via process reward guided tree search. *Advances in Neural Information Processing Systems*, 37:64735–64772.
- Kaiyan Zhang, Yuxin Zuo, Bingxiang He, Youbang Sun, Runze Liu, Che Jiang, Yuchen Fan, Kai Tian, Guoli Jia, Pengfei Li, et al. 2025. A survey of reinforcement learning for large reasoning models. *arXiv preprint arXiv:2509.08827*.

## A Appendix

### A.1 Generalization Analysis

To further investigate the gap, we analyze the relationship between training and validation accuracy on MedQA for Qwen3-4B (Figure 9). We find that PPO exhibits clear overfitting: its training accuracy continues to increase while validation accuracy plateaus, indicating weak generalization. In contrast, both GRPO and TEMPO improve training and validation accuracy in tandem, with TEMPO achieving the highest performance on both, suggesting more reliable generalization rather than memorization of the training distribution. This analysis explains why PPO lags behind in final MedQA accuracy and underscores TEMPO’s advantage when scaling to larger models and more challenging domains.

### A.2 Datasets and Models

We consider two domains: mathematics and medicine. For training, we adopt one representative dataset from each domain: MATH (Hendrycks et al.) for mathematical reasoning and MedQA (Jin et al., 2021) for medical question answering. These serve as the *in-distribution* (ID) training tasks. For evaluation, we test on both the in-distribution test sets of MATH and MedQA, as well as multiple *out-of-distribution* (OOD) benchmarks to assess generalization. In the math domain, OOD benchmarks include GSM-HARD (Gao et al., 2023), a challenging variant of GSM8K with harder grade-school problems, and AMC23<sup>2</sup>, a set of recent American Mathematics Competition problems. In the medical domain, OOD benchmarks include MedMCQA (Pal et al., 2022), a dataset consisting of multiple-choice medical questions designed to test clinical knowledge, and MMLU-Medical (Singhal et al., 2023), a medical subset of the Massive Multi-task Language Understanding (MMLU) benchmark focusing on diverse topics in the medical field. We adopt two publicly available models from the Qwen 3 (Yang et al., 2025a) family: Qwen3-1.7B and Qwen3-4B. Both models are fine-tuned in our experiments under identical settings to ensure fair comparison.

### A.3 Implementation Details

To ensure our GRPO implementation is robust, and our evaluation reflects its full potential, we have

<sup>2</sup><https://huggingface.co/datasets/AI-MO/aimo-validation-amc>

applied a set of well-established techniques and best practices from the literature (Yu et al., 2025). Below, we outline the key implementation details that were most effective in our experiments:

- **Clip-Higher** (decoupled clipping). We decouple the clipping bounds and raise the upper cap ( $1 + \epsilon_{\text{high}}$ ) while keeping the lower cap ( $1 - \epsilon_{\text{low}}$ ), which allows low-probability “exploration” tokens to increase more freely and helps prevent entropy collapse.
- **Token-level policy-gradient loss**: Token-level policy-gradient loss (global token averaging): We optimize a token-level surrogate averaged over *all* tokens in the batch, broadcasting each response’s group-normalized outcome reward to its tokens since sample-level averaging underweights long responses and fails to penalize low-quality long patterns, which destabilizes training; token-level loss restores balanced credit assignment and yields healthier length/entropy dynamics.
- **Remove KL divergence**: In long-CoT reasoning, the online policy can beneficially diverge from the initialization; thus we omit an explicit KL regularizer and rely on clipping for stability.

**Training Details and Hyperparameters** We adopt a binary task reward  $R$  that evaluates final answer correctness against ground truth, following previous work (Huang et al., 2024; Ivison et al., 2024). To ensure fair comparison, all methods consume the same number of episodes during training: for each question, we sample 6 episodes and go over the dataset 10 times, yielding 60 episodes per question across all methods.

### A.4 Hyperparameter

In this section, we provide a comprehensive overview of the hyperparameters used in our experiments. The number of training episodes was carefully selected to ensure that the amount of training data remained consistent across all methods.

**PPO** Finetuning LLMs with PPO is known to be highly sensitive to hyperparameter choices, making optimal selection critical for strong performance. To ensure robustness, we considered hyperparameter values reported in prior studies (Shao et al., 2024) and performed extensive sweeps across a wide range of candidate values. Specifically,

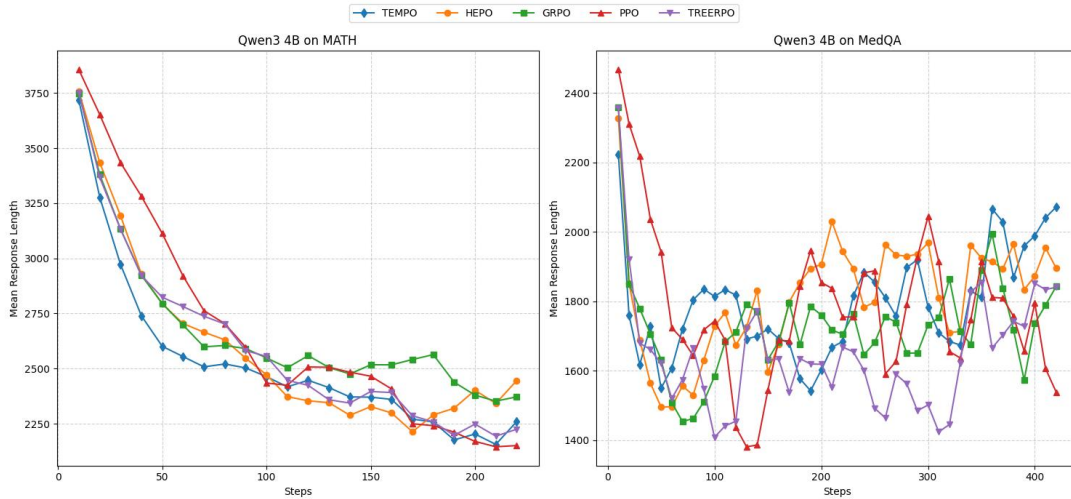


Figure 8: Mean Responses Length of MATH and MedQA for Qwen3- Qwen3-4B accross training.

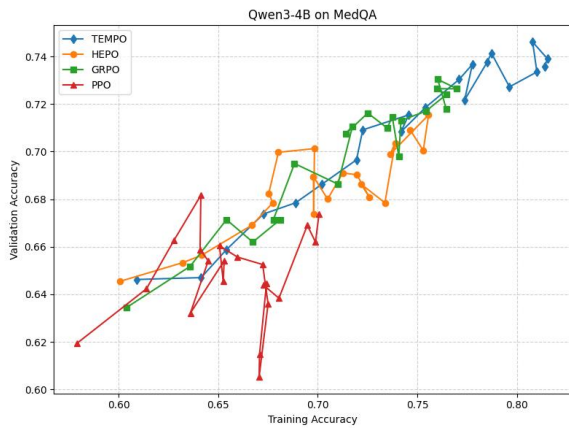


Figure 9: Training vs. validation accuracy on MedQA with Qwen3-4B. PPO overfits to training data, while TEMPO maintains better generalization.

we first identified the set of hyperparameters that achieved the best performance across both the MATH and MedQA tasks using the Qwen3 1.7B model. This optimal configuration was then employed for the remainder of our experiments. The complete list of PPO hyperparameters, along with their respective search spaces, is shown in Table 2.

**GRPO, HEPO, and TEMPO** Since policy optimization in RLOO and GRPO closely resembles PPO, we initialized their hyperparameters using the PPO configuration. This ensures a strong starting point while enabling a systematic comparison among the algorithms. We note that the absence of explicit credit assignment in these methods may result in high-variance policy gradient updates, potentially leading to instability (Greensmith et al., 2004). The full list of hyperparameters for GRPO,

HEPO, and TEMPO is provided in Table 2.

### A.5 Compute

All experiments were conducted using multi-GPU training to efficiently handle the computational demands of large-scale models. For the Qwen3-1.7B model, we utilized a node with 1 × Nvidia H100 80GB GPUs to train both TEMPO and all the baselines. For the larger Qwen3-4B model, we employed a more powerful setup, using a node with 2 × Nvidia H100 80GB GPUs.

### A.6 Software Stack

For model implementation, we utilize the Hugging-face library. Training is carried out using the VERL (Zhang et al., 2024a) distributed training library, which offers efficient multi-GPU support. For trajectory sampling during RL training, we rely on the vLLM library (Kwon et al., 2023), which provides optimized inference for LLMs.

### A.7 Reproducibility

In this study, all experiments were conducted using open-source libraries, publicly available datasets, and open-weight LLMs. To ensure full reproducibility, we will make our codebase publicly available on GitHub at <https://github.com/fatebreaker/tempo>.

### A.8 LLM Usage

In accordance with the ICLR 2026 policies on LLM usage, we disclose how LLMs were used in this work. LLMs were employed to assist with grammar polishing, wording improvements, and drafting text during paper preparation. All technical content,

proofs, experiments, and analyses were conceived, implemented, and validated by the authors. Authors remain fully responsible for the correctness of the claims and results.

No LLMs were used to generate research ideas, write code for experiments, or produce results. No confidential information was shared with LLMs, and no prompt injections or other inappropriate uses were involved.

This disclosure aligns with the ICLR Code of Ethics: contributions of tools are acknowledged, while accountability and verification rest entirely with the human authors.

Parameter	Value	Notes
<b>Training</b>		
Optimizer	AdamW	
Adam parameters $(\beta_1, \beta_2)$	(0.9, 0.999)	
Learning rate	$1 \times 10^{-6}$	
Weight decay	0.0	
Warmup	0% of training steps	
# Train steps (MATH)	220 steps	~10 dataset epochs
# Train steps (MedQA)	420 steps	~10 dataset epochs
<b>General</b>		
Maximum prompt length	1024 tokens	
Maximum response length	8192 tokens	
Training batch size	512	
<b>PPO</b>		
Mini-batch size	64	
# Inner epochs per PPO step	2	
Discount factor $\gamma$	1.0	
GAE parameter $\lambda$	1.0	
KL penalty coefficient $\beta$	$1 \times 10^{-4}$	
<b>GRPO/HEPO/TEMPO</b>		
# Responses per prompt	6	
Mini-batch size	64	
Discount factor $\gamma$	1.0	
KL penalty coefficient $\beta$	0.0	
Policy clipping parameter $\epsilon$	0.28, 0.2	
<b>HEPO</b>		
$\rho$	0.2	Only do gradient update on top 20% high entropy tokens

Table 2: Summary of hyperparameters used in the experiments.