# OptPO: Optimal Rollout Allocation for Test-time Policy Optimization

**Youkang Wang** [1]  **Jian Wang** [1]  **Rubing Chen** [1]  **Tianyi Zeng** [1]  **Xiao-Yong Wei** [2 1]  **Qing Li** [1]

## Abstract

Test-time policy optimization enables large language models (LLMs) to adapt to distribution shifts by leveraging feedback from self-generated rollouts. However, existing methods rely on fixed-budget majority voting to estimate rewards, incurring substantial computational redundancy. We propose Optimal Rollout Allocation for Test-time Policy Optimization (OptPO), a principled framework that adaptively allocates inference budgets. By formulating the voting process as a Bayesian sequential probability ratio test, OptPO dynamically halts sampling once the posterior confidence in a consensus answer exceeds a specified threshold. Crucially, it utilizes the retained rollouts for on-policy updates, seamlessly integrating with algorithms like PPO or GRPO without requiring ground-truth labels. Across diverse reasoning benchmarks, OptPO significantly reduces rollout overhead compared to fixed-sample baselines while preserving or improving accuracy. By unifying statistically optimal stopping with test-time learning, OptPO offers a computationally efficient paradigm for test-time adaptation. The source code will be open upon acceptance at https://open-upon-acceptance.

## 1. Introduction

Large language models (LLMs) have achieved remarkable advances in natural language understanding, reasoning, and problem solving (Brown et al., 2020; Achiam et al., 2023). Yet, a persistent challenge emerges when models face distribution shifts at inference time, such as domain-specific terminology, unfamiliar reasoning styles, or linguistic variability, which can cause substantial performance degradation (Hu et al., 2025). To mitigate this, test-time adaptation that addresses these shifts dynamically has thus emerged as a critical research direction.

As a promising paradigm, *test-time learning* (TTL) enables models to learn to adapt during inference. Early methods primarily focused on light-weight adaptation without updating parameters, such as entropy minimization or batch normalization (Wang et al., 2020; Niu et al., 2022). More recent work extends these ideas to *test-time training* (TTT), with explicit parameter updates occurring at inference, which substantially enhances task-specific adaptability. For example, Akyürek et al. (2024) suggests that test-time training improves few-shot learning under distribution shift, while Hübotter et al. (2024) proposes active fine-tuning that selectively adapts LLMs using informative feedback signals.

Within the TTT paradigm, a particularly emerging direction is test-time policy optimization. By treating the inference process as a reinforcement learning (RL) problem, methods such as RL with verifiable rewards (RLVR) (Wang et al., 2025) and TTRL (Zuo et al., 2025) utilize consensus-based feedback (e.g., majority voting) to guide policy updates. These approaches leverage the "wisdom of the crowd" from multiple rollouts to construct label-free rewards, effectively bridging the gap between reasoning consistency and correctness. However, a critical bottleneck remains: *rollout inefficiency*. Standard test-time RL methods typically mandate a fixed, large number of rollouts for every query to estimate rewards robustly. This "one-size-fits-all" strategy is computationally wasteful, as easy instances establish consensus quickly, while only hard instances require extensive sampling. Consequently, the heavy inference cost limits the practical deployment of test-time policy optimization.

This paper introduces *Optimal Rollout Allocation for Test-time Policy Optimization* (**OptPO**), a general framework designed to address the efficiency-accuracy trade-off. Our key insight is to reframe reward estimation as a sequential hypothesis testing problem. Instead of pre-committing to a fixed sample budget, OptPO adaptively estimates the vote distribution using a Bayesian Sequential Probability Ratio Test (SPRT). Sampling halts as soon as the posterior evidence for a consensus answer satisfies a user-defined error tolerance. The accumulated rollouts are then immediately repurposed for policy optimization (e.g., via PPO (Schulman et al., 2017) or GRPO (Shao et al., 2024)), ensuring that compute is allocated only where necessary to reduce uncertainty. Furthermore, our OptPO can be broadly compatible with supervised fine-tuning for test-time training.

---

[1]Department of Computing, The Hong Kong Polytechnic University [2]Department of Computer Science, Sichuan University. Correspondence to: Xiao-Yong Wei <cs007.wei@polyu.edu.hk>.

In summary, our contributions are as follows:

- We identify the fixed-budget rollout mechanism as a primary source of inefficiency in current test-time policy optimization approaches.

- We propose OptPO, a principled, plug-and-play framework for test-time policy optimization that adaptively allocates rollouts based on sequential majority voting with early stopping.

- We demonstrate that OptPO substantially reduces computational costs (e.g., over 40% savings in token consumption on the GPQA benchmark), while maintaining competitive accuracy across representative reasoning benchmarks.

By bridging ideas from test-time learning and reinforcement learning, we advance the broader goal of enabling LLMs to self-improve efficiently in dynamic environments.

## 2. Related Work

**Test-Time Adaptation for LLMs.** Traditional machine learning assumes a fixed training–testing separation, with frozen model parameters at inference. However, models often encounter distribution shifts at test time, leading to degraded performance. Early approaches to *test-time adaptation* (TTA) proposed lightweight adjustments such as entropy minimization or updating normalization statistics on unlabeled test batches (Wang et al., 2020; Niu et al., 2022). These methods demonstrated effectiveness on computer vision benchmarks but remain limited in expressiveness for large language models (LLMs), where domain shift can involve complex reasoning and linguistic variability.

**Test-Time Learning and Training.** Recent work extends adaptation to full or partial parameter updates during inference. Akyürek et al. (2024) shows that test-time training (TTT) substantially improves the few-shot generalization of LLMs under distribution shift. Hübotter et al. (2024) propose active fine-tuning strategies that selectively adapt LLMs using informative examples. Retrieval-based test-time training offers another perspective: Hardt & Sun (2023) demonstrates that fine-tuning on nearest neighbors retrieved from a large corpus before answering each test instance (TTT-NN) significantly boosts language modeling performance. Similarly, Hu et al. (2025) proposes perplexity-based test-time adaptation for LLMs, framing inference as a continual self-supervised training process. Collectively, these works illustrate the broad potential of TTT in improving the adaptability of large models.

**Test-Time Reinforcement Learning.** Reinforcement learning has also been employed for test-time training, with the incorporation of verification or feedback signals. Li et al. (2025) introduces feedback-guided test-time training, where external verifiers provide weak supervision to refine predictions during inference. Moreover, Zuo et al. (2025) presents Test-Time Reinforcement Learning (TTRL), which generates multiple rollouts, aggregates predictions through majority voting, and uses the consensus as a reward signal to update the policy. While effective, these methods suffer from inefficiency: they require many redundant rollouts to stabilize consensus estimates. More recently, Wang et al. (2025) shows that reinforcement learning with verifiable rewards (RLVR) can achieve strong reasoning gains with extremely few training examples, underscoring the potential of efficient RL in test-time settings. Our work builds on this line by explicitly addressing rollout inefficiency, introducing a statistically grounded early-stopping mechanism for adaptive allocation.

**Policy Optimization in RL.** Policy optimization techniques have been widely used in LLM-based reinforcement learning. Proximal Policy Optimization (PPO) (Schulman et al., 2017) is a standard algorithm for stabilizing policy updates under KL constraints and has been widely applied in reinforcement learning from human feedback. Group Relative Policy Optimization (GRPO) (Shao et al., 2024) extends this paradigm with group-normalized advantages, enabling more stable optimization for LLMs under limited feedback. OptPO leverages the generality of these surrogates, providing a plug-and-play rollout allocation mechanism that enhances efficiency without requiring changes to the underlying optimization algorithm.

In summary, prior work has established test-time training as a powerful framework for improving model adaptability during inference, spanning entropy minimization, verifier-guided learning, and reinforcement learning. However, existing test-time policy optimization methods remain computationally inefficient due to excessive rollout requirements. Our work fills this gap by introducing OptPO, a framework for optimal rollout allocation that substantially reduces compute costs while preserving performance, complementing and extending the existing test-time adaptation methods.

## 3. Methodology

### 3.1. Preliminaries

We consider test-time training for tasks with deterministic, verifiable targets (e.g., mathematics or science reasoning problems). Each test instance provides a prompt $x \in \mathcal{X}$, and the model $\pi_\theta(\cdot \mid x)$, parameterized by $\theta$, defines a stochastic policy over completions $y \in \mathcal{Y}$. Although ground-truth labels are unavailable at test time, we can estimate a self-consistent pseudo-label from multiple rollouts to guide adaptation. The objective of test-time training is divided into

two steps: **(i) Pseudo-label estimation:** efficiently identify a consensus $a^\star$ via early-stopped majority voting; **(ii) Parameter update:** adapt $\pi_\theta$ using self-supervised feedback defined by $a^\star$, via either RL- or SFT-based updates.

**Rollouts and Votes.** Each completion $y$ is mapped to a canonical answer token $a = \mathcal{E}(y)$ through an *answer extractor* $\mathcal{E} : \mathcal{Y} \to \mathcal{A}$, where $\mathcal{A} = \{1, \ldots, m\}$ denotes the finite set of possible answers. We model $m$ candidate hypotheses $\{H_1, \ldots, H_m\}$, where $H_j$ represents the hypothesis that the correct answer equals $j$. Sequentially generating completions yields a sequence of *votes*

$$D_t = \mathcal{E}(Y_t), \qquad Y_t \sim \pi_\theta(\cdot \mid x),$$

with $D_t \in \{1, \ldots, m\}$ and cumulative vote history $D_{1:t} = \{D_1, \ldots, D_t\}$. Each instance is allowed a maximum rollout budget $M$ and a minimum retained sample count $N \leq M$. The goal is to form a reliable consensus pseudo-label $a^\star$ using as few rollouts as possible, then update $\theta$ at test time using a small set of $N$ retained rollouts.

**Answer Probabilities.** Conditioned on the true hypothesis $H_j$, each vote equals $j$ (the correct answer) with probability $p_0 \in (0, 1)$, and one of the remaining $m-1$ candidates otherwise, with probabilities $\{p_r\}_{r=1}^{m-1}$ satisfying $p_0 + \sum_r p_r = 1$ and $p_0 > p_r$. Votes are conditionally independent given $H_j$. This categorical noise model captures the intuition that correct answers appear more often and provides a one-dimensional sufficient statistic for sequential testing based on vote-count gaps.

## 3.2. Rollout Allocation with Probabilistic Optimality

We treat each answer candidate as a hypothesis and perform a *sequential probability ratio test* (SPRT) between the current leading and second-place classes as rollouts accumulate.

**Likelihood and Posterior.** Let $v_j(t)$ denote the number of votes for class $j$ up to time $t$, and let $L(t)$ and $S(t)$ denote the leading and runner-up classes. Under hypothesis $H_j$, the likelihood of the observed votes factorizes as

$$P(D_{1:t} \mid H_j) = \prod_{k=1}^{t} P(D_k \mid H_j) = p_0^{v_j(t)} \prod_{c \neq j} p_c^{v_c(t)}.$$

With uniform priors $P(H_j) = 1/m$, the posterior satisfies $P(H_j \mid D_{1:t}) \propto P(D_{1:t} \mid H_j)$.

**Bayes Likelihood Ratio between Top Two Classes.** We define the Bayes factor (likelihood ratio under uniform priors) as

$$\Lambda_t = \frac{P(H_{L(t)} \mid D_{1:t})}{P(H_{S(t)} \mid D_{1:t})} = \frac{P(D_{1:t} \mid H_{L(t)})}{P(D_{1:t} \mid H_{S(t)})}.$$

Expanding and simplifying the equation yields:

$$\Lambda_t = \frac{p_0^{v_L(t)} \prod_{c \neq L} p_c^{v_c(t)}}{p_0^{v_S(t)} \prod_{c \neq S} p_c^{v_c(t)}} = \frac{p_0^{v_L(t)}}{p_0^{v_S(t)}} \cdot \frac{\left(p_{S|H_L}\right)^{v_S(t)}}{\left(p_{L|H_S}\right)^{v_L(t)}}.$$

Crucially, the per-class error probabilities in the numerator and denominator are *conditioned on different hypotheses*. Under $H_{L(t)}$, every incorrect class (including $S(t)$) shares the same probability mass. And symmetrically, under $H_{S(t)}$, the (now incorrect) leader $L(t)$ is one of the $m-1$ wrong classes, hence

$$p_{S|H_L} = \frac{1 - p_0}{m - 1}, \qquad p_{L|H_S} = \frac{1 - p_0}{m - 1}.$$

Therefore $p_{S|H_L} = p_{L|H_S} = (1 - p_0)/(m - 1)$, and the extra factor simplifies to:

$$\frac{\left(p_{S|H_L}\right)^{v_S(t)}}{\left(p_{L|H_S}\right)^{v_L(t)}} = \left(\frac{1 - p_0}{m - 1}\right)^{v_S(t) - v_L(t)}$$

$$= \left(\frac{m - 1}{1 - p_0}\right)^{v_L(t) - v_S(t)}.$$

Combining terms, the Bayes factor is derived as:

$$\Lambda_t = \left[\frac{p_0(m - 1)}{1 - p_0}\right]^{\Delta_t} \equiv \kappa^{\Delta_t}, \qquad \kappa = \frac{p_0(m - 1)}{1 - p_0},$$

where $\Delta_t = v_L(t) - v_S(t)$ is the vote-count gap.

**Stopping Criterion.** Given user-specified type-I and type-II error budgets $\alpha, \beta \in (0, 1)$, we define Wald thresholds as follows:

$$A = \frac{1 - \beta}{\alpha}, \qquad B = \frac{\beta}{1 - \alpha}, \qquad (A > 1 > B).$$

At each step $t$, stop and select $L(t)$ if $\Lambda_t \geq A$, select $S(t)$ if $\Lambda_t \leq B$, and otherwise continue (subject to the hard cap $M$). Since $\Lambda_t = \kappa^{\Delta_t}$ with $\Delta_t \in \mathbb{Z}_{\geq 0}$, we can equivalently threshold on the integer gap $\Delta_t$:

$$\Delta_A = \left\lceil \frac{\log A}{\log \kappa} \right\rceil, \qquad \Delta_B = \left\lfloor \frac{\log B}{\log \kappa} \right\rfloor.$$

If $\kappa > 1$ (typical when $p_0 > \frac{1}{m}$), then $\Lambda_t$ increases with $\Delta_t$: stop for $L$ when $\Delta_t \geq \Delta_A$ and for $S$ when $\Delta_t \leq \Delta_B$. If $\kappa < 1$ (pathological low-accuracy regime), the inequalities reverse; in practice we ensure $p_0 > \frac{1}{m}$ so $\kappa > 1$.

We denote the stopping time by

$$\tau = \min\left\{ N \leq t \leq M : \Delta_t \notin (\Delta_B, \Delta_A) \right\}.$$

The final consensus pseudo-label is $a^\star = L(\tau)$ if $\Lambda_\tau \geq A$, or $a^\star = S(\tau)$ if $\Lambda_\tau \leq B$; otherwise if no boundary is crossed by $M$, we set $a^\star = L(M)$.

This procedure adaptively determines how many rollouts each instance requires to reach a confident majority, forming the backbone of our efficient test-time training framework.

**Rollout Allocation.** Sampling proceeds until $t = \max\{N, \tau\}$, guaranteeing at least $N$ retained rollouts for learning while skipping redundant ones on easier instances. Saved compute can be reallocated to harder examples under a fixed global budget.

## 3.3. Test-Time Policy Optimization

Once a confident pseudo-label $a^\star$ is obtained, test-time adaptation updates the parameters $\theta$ based on the $N$ retained rollouts. The adaptation rule differentiates two paradigms depending on the update formulation.

**RL-Based Optimization.** For reinforcement-style updates, each retained rollout $(Y_i, a_i)$ receives a label-free reward measuring alignment with the consensus:

$$R(Y_i; a^\star) = \mathbb{I}[a_i = a^\star],$$

or more generally, $R(Y_i; a^\star) \in [0, 1]$ for partial credit or verification heuristics. The test-time policy optimization objective for prompt $x$ is

$$\max_\theta \ \mathbb{E}_{Y \sim \pi_\theta(\cdot|x)}[R(Y; a^\star)].$$

With policy-gradient updates and baseline $b_i$, one on-policy step over $N$ samples is

$$\theta \leftarrow \theta + \eta \frac{1}{N} \sum_{i=1}^{N} \nabla_\theta \log \pi_\theta(Y_i \mid x) \left( R(Y_i; a^\star) - b_i \right),$$

optionally regularized to a reference policy $\pi_{\mathrm{ref}}$:

$$\theta \leftarrow \arg\max_\theta \frac{1}{N} \sum_{i=1}^{N} \Big[ \log \pi_\theta(Y_i \mid x)\, \hat{A}_i \\ - \beta_{\mathrm{KL}}\, \mathrm{KL}\big(\pi_\theta(\cdot \mid x) \,\|\, \pi_{\mathrm{ref}}(\cdot \mid x)\big) \Big], \quad (1)$$

where $\hat{A}_i = R(Y_i; a^\star) - b_i$. This objective is compatible with PPO (Schulman et al., 2017), GRPO (Shao et al., 2024), or any general RL surrogates. The OptPO algorithm thus combines (i) BSPRT-based early stopping for efficient rollout allocation and (ii) on-policy policy optimization for rapid test-time adaptation.

**SFT-Based Optimization.** In the supervised variant, parameter updates are performed directly toward the pseudo-labels. For each batch of test inputs $\{x_i\}_{i=1}^{B}$, the model generates rollouts and determines a consensus label $G_i$ via majority voting. The standard test-time supervised fine-tuning (TTSFT) objective is:

$$\mathcal{L}_{\mathrm{TTSFT}}(\theta) = -\frac{1}{B} \sum_{i=1}^{B} \log \pi_\theta(G_i \mid x_i).$$

To improve efficiency, our *OptPO-SFT* improves the fixed-$N$ majority voting with the BSPRT-based early-stopping rule described earlier, producing adaptive pseudo-labels $G_i = L(\tau_i)$ where $\tau_i$ is determined by the likelihood ratio thresholds $(A, B)$. The same cross-entropy objective applies:

$$\mathcal{L}_{\mathrm{OptPO-SFT}}(\theta) = -\frac{1}{B} \sum_{i=1}^{B} \log \pi_\theta(G_i \mid x_i).$$

Thus, OptPO-SFT shares the same update rule as TTSFT but allocates rollouts adaptively via the Bayesian stopping criterion, unifying efficiency with reliability.

**Unified Perspective.** Both RL-based and SFT-based variants of OptPO instantiate the same three-stage framework: They initiate *rollouts*, compute *BSPRT-based pseudo-label*, and finally perform *parameter update*. They differ by the update rule: policy-gradient optimization versus cross-entropy fine-tuning, making the framework conceptually unified.

## 4. Experiments

### 4.1. Experimental Setup

**Benchmarks.** We adopt challenging math and science reasoning benchmarks: (1) **MATH-500** (Hendrycks et al., 2021), with 500 high school competition problems; (2) **AIME 2024** (MAA, 2024), containing 30 pre-Olympiad level problems; (3) **AMC**, with intermediate-level problems from American math competitions. (4) **GPQA** (Rein et al., 2024), a Q&A dataset of challenging, graduate-level science questions in biology, physics, and chemistry.

**Backbone Models.** We evaluate OptPO on widely used open-source LLM backbones, including Qwen2.5-Math-1.5B, Qwen2.5-7B, and Llama-3.2-1B-Instruct, to cover various model sizes and verify cross-architecture generalization. When verifying the effectiveness of TTSFT and OptPO-SFT, we utilized Qwen2.5-Math-1.5B and Llama-3.1-8B-Instruct as backbone models.

**Baseline Methods.** OptPO is a plug-and-play framework that can be plugged into multiple existing RL frameworks for improving test-time training efficiency. We plugged it into **PPO** (Schulman et al., 2017), **GRPO** (Shao et al., 2024), and Reinforce++ (**Rei++**) (Hu, 2025) for the comprehensiveness of evaluation. We compare its performance with TTRL (Zuo et al., 2025), the most prevalent test-time policy optimization method. We also evaluated the performance of OptPO-SFT in similar settings.

**Implementation Details.** Each rollout uses temperature 0.6 and nucleus sampling ($p = 0.95$). We generate up to

*Table 1.* Performance comparisons of OptPO with existing test-time reinforcement learning methods.

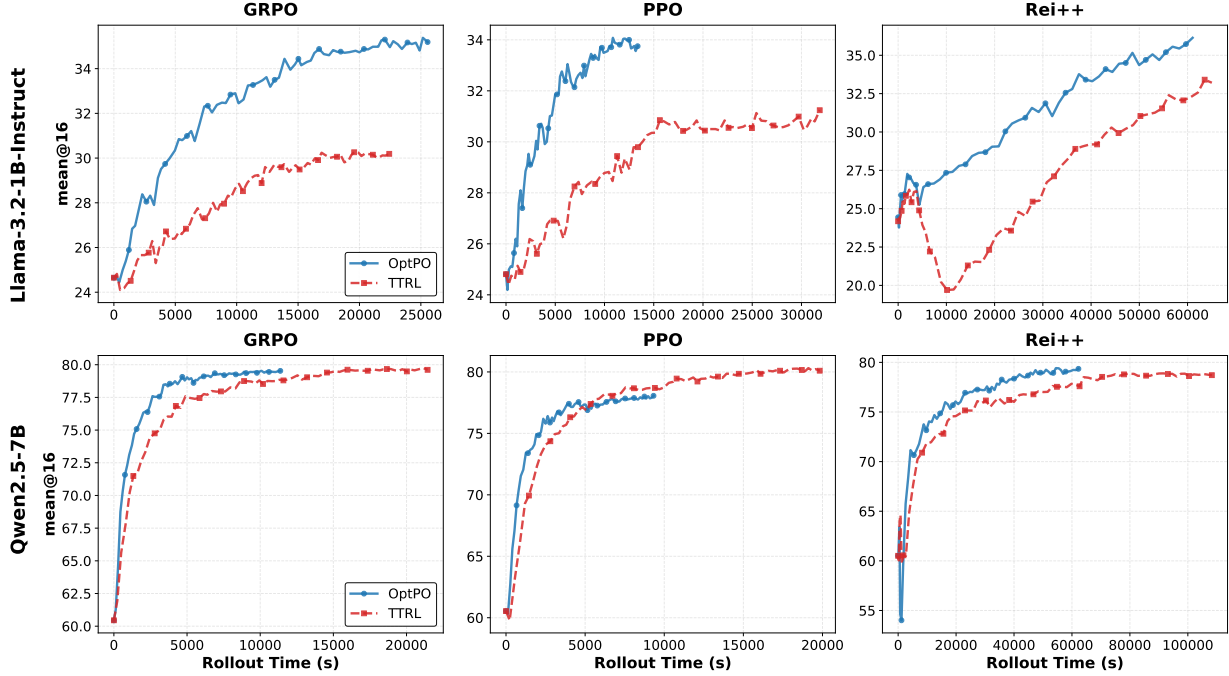| Benchmark | AIME 2024 | | | | | AMC | | | | | MATH-500 | | | | | GPQA | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | Mean @16 | Pass @16 | Pass @1 | Toks. | Toks. saving | Mean @16 | Pass @16 | Pass @1 | Toks. | Toks. saving | Mean @16 | Pass @16 | Pass @1 | Toks. | Toks. saving | Mean @16 | Pass @16 | Pass @1 | Toks. | Toks. saving |
| **Qwen2.5-Math-1.5B** | | | | | | | | | | | | | | | | | | | | |
| TTRL-GRPO | 15.6 | 43.3 | 20.0 | 142M | 17.25% | 47.7 | 83.1 | 54.2 | 159M | 26.93% | 72.3 | 90.6 | 73.6 | 218M | 32.57% | 26.5 | 87.8 | 30.5 | 251M | 43.20% |
| OptPO-GRPO | 14.6 | 46.7 | 20.0 | 117M | | 48.5 | 81.9 | 54.2 | 116M | | 72.3 | 90.8 | 73.8 | 147M | | 26.8 | 87.8 | 34.5 | 143M | |
| TTRL-PPO | 18.5 | 46.7 | 23.3 | 143M | 15.20% | 48.1 | 83.1 | 55.4 | 160M | 26.80% | 71.9 | 91.2 | 73.8 | 231M | 35.01% | 25.9 | 87.8 | 29.9 | 262M | 45.06% |
| OptPO-PPO | 16.0 | 46.7 | 23.3 | 121M | | 48.1 | 83.1 | 53.0 | 117M | | 72.3 | 91.0 | 73.8 | 150M | | 26.6 | 89.8 | 32.5 | 144M | |
| TTRL-Rei++ | 11.9 | 43.3 | 20.0 | 198M | 14.35% | 39.3 | 81.9 | 49.4 | 242M | 27.42% | 55.9 | 90.2 | 56.4 | 696M | 25.53% | 26.3 | 93.9 | 32.0 | 431M | 43.52% |
| OptPO-Rei++ | 10.4 | 46.7 | 20.0 | 170M | | 38.3 | 81.9 | 44.6 | 176M | | 56.2 | 89.8 | 57.6 | 518M | | 26.5 | 90.9 | 28.9 | 243M | |
| **Llama-3.2-1B-Instruct** | | | | | | | | | | | | | | | | | | | | |
| TTRL-GRPO | 3.5 | 23.3 | 6.7 | 134M | 47.62% | 14.9 | 51.8 | 18.1 | 128M | 78.36% | 30.3 | 66.0 | 30.8 | 181M | 14.77% | 26.4 | 82.2 | 28.4 | 134M | 44.74% |
| OptPO-GRPO | 1.7 | 20.0 | 6.7 | 70M | | 15.2 | 55.4 | 16.9 | 28M | | 35.4 | 64.4 | 36.0 | 155M | | 25.2 | 81.2 | 26.9 | 74M | |
| TTRL-PPO | 3.8 | 30.0 | 6.7 | 47M | 13.70% | 13.9 | 55.4 | 19.3 | 75M | 41.76% | 31.2 | 64.6 | 31.6 | 246M | 44.49% | 25.7 | 83.2 | 28.9 | 159M | 46.19% |
| OptPO-PPO | 3.8 | 30.0 | 6.7 | 41M | | 17.1 | 54.2 | 22.9 | 44M | | 34.1 | 64.4 | 34.8 | 137M | | 26.7 | 81.7 | 31.0 | 86M | |
| TTRL-Rei++ | 1.5 | 16.7 | 6.7 | 44M | 30.02% | 11.4 | 54.2 | 16.9 | 32M | 18.96% | 33.4 | 65.4 | 32.2 | 557M | 29.40% | 27.1 | 85.8 | 30.5 | 478M | 41.96% |
| OptPO-Rei++ | 1.7 | 23.3 | 6.7 | 31M | | 11.0 | 51.8 | 13.3 | 26M | | 36.1 | 66.6 | 38.0 | 393M | | 26.7 | 85.8 | 28.9 | 278M | |
| **Qwen2.5-7B** | | | | | | | | | | | | | | | | | | | | |
| TTRL-GRPO | 23.5 | 40.0 | 23.3 | 87M | 38.56% | 52.1 | 78.3 | 54.2 | 111M | 49.06% | 79.7 | 91.0 | 80.0 | 172M | 45.02% | 35.1 | 85.3 | 39.1 | 98M | 40.31% |
| OptPO-GRPO | 23.5 | 43.3 | 23.3 | 53M | | 56.8 | 77.1 | 57.8 | 56M | | 79.5 | 91.6 | 79.8 | 94M | | 35.7 | 84.8 | 38.6 | 59M | |
| TTRL-PPO | 24.4 | 46.7 | 23.3 | 104M | 42.85% | 54.3 | 79.5 | 55.4 | 101M | 41.66% | 80.3 | 91.0 | 80.4 | 171M | 48.96% | 37.9 | 83.8 | 40.1 | 109M | 39.83% |
| OptPO-PPO | 24.0 | 43.3 | 26.7 | 60M | | 53.2 | 79.5 | 55.4 | 59M | | 78.1 | 91.0 | 78.4 | 87M | | 36.8 | 88.8 | 39.1 | 66M | |
| TTRL-Rei++ | 20.8 | 40.0 | 23.3 | 345M | 41.45% | 52.9 | 80.7 | 55.4 | 430M | 43.50% | 78.9 | 90.8 | 80.0 | 886M | 44.52% | 33.0 | 85.3 | 32.5 | 139M | 45.47% |
| OptPO+Rei++ | 20.6 | 40.0 | 23.3 | 202M | | 56.0 | 79.5 | 57.8 | 243M | | 79.4 | 91.0 | 80.6 | 491M | | 33.0 | 84.3 | 34.0 | 76M | |



*Figure 1.* Performance (mean@16 accuracy vs rollout time) comparison of OptPO and TTRL on the MATH-500 benchmark, integrating different RL algorithms and using (i) Llama-3.2-1B-Instruct (top row) and (ii) Qwen2.5-7B (bottom row) as backbones.

$M = 64$ rollouts per instance with a minimum of $N = 32$ or $N = 16$ samples used for policy update. For RL-based updates (OptPO), we perform one policy-gradient step per instance using PPO/GRPO-style objectives with a KL penalty $\beta_{KL} = 0.001$. Actor Learning Rate is set to 1e-6, for RL algorithms involving a critic (like PPO), the critic learning rate is set to 1e-5. For SFT-based training (including TTSFT and OptPO-SFT), gradient descent is applied on the pseudo-labeled batch using AdamW with a learning rate of 2e-5.

The probability of answering a specific question $x$ correctly $p_0$ is estimated adaptively as the ratio of majority answer

*Table 2.* Performance comparisons between OptPO-SFT and existing test-time supervised fine-tuning (TTSFT).

| Benchmark | AIME 2024 | | | | | AMC | | | | | MATH-500 | | | | | GPQA | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | mean @16 | pass @16 | pass @1 | toks. | toks. saving | mean @16 | pass @16 | pass @1 | toks. | toks. saving | mean @16 | pass @16 | pass @1 | toks. | toks. saving | mean @16 | pass @16 | pass @1 | toks. | toks. saving |
| **Qwen2.5-Math-1.5B** | | | | | | | | | | | | | | | | | | | | |
| TTSFT | 13.1 | 46.7 | 20.0 | 166M | 0.61% | 37.9 | 80.7 | 42.2 | 138M | 12.08% | 59.7 | 90.8 | 63.8 | 200M | 25.29% | 27.9 | 81.8 | 30.3 | 113M | 51.17% |
| OptPO-SFT | 13.5 | 50.0 | 20.0 | 165M | | 39.7 | 81.9 | 43.4 | 121M | | 59.0 | 91.2 | 60.2 | 150M | | 27.6 | 78.8 | 31.3 | 55M | |
| **Llama-3.1-8B-Instruct** | | | | | | | | | | | | | | | | | | | | |
| TTSFT | 3.8 | 26.7 | 6.7 | 314M | 15.33% | 20.1 | 63.9 | 28.9 | 326M | 17.16% | 43.5 | 81.2 | 44.2 | 653M | 19.84% | 28.8 | 85.9 | 30.3 | 382M | 61.89% |
| OptPO-SFT | 6.2 | 36.7 | 3.3 | 266M | | 18.7 | 60.2 | 19.3 | 270M | | 43.6 | 82.8 | 47.0 | 523M | | 28.9 | 82.8 | 32.3 | 146M | |

among the minimum retained rollouts, then times a degradation factor 0.6 to obtain a more realistic estimate of $p_0$, since instinctively pseudo-accuracy estimated by majority probability is higher than real accuracy. This calibration stabilizes BSPRT thresholds without labeled data. We set $(\alpha, \beta) = (0.05, 0.05)$, yielding a 95% confidence level in the consensus decision. We require exceeding this confidence level 5 times across the rollouts for rigorous guarantees and robust early-stopping. All experiments are conducted on $8 \times$ H20 GPUs.

## 4.2. Main Results

Table 1 reports the performance of our proposed methods across benchmarks and backbones. Overall, OptPO consistently outperforms its non-optimized counterpart TTRL, while substantially reducing rollout costs, demonstrating the effectiveness of our Bayesian early-stopping and adaptive pseudo-labeling mechanisms. Figure 1 shows faster convergence and improvement of OptPO on the MATH-500 benchmark when using different backbone models. Our key findings are as follows:

**Substantial efficiency gains with comparable or superior accuracy.** OptPO achieves comparable or higher mean@16 and pass@16 scores than baseline TTRL methods across all benchmarks, while consuming 30% $\sim$ 50% fewer tokens on average. For instance, on GPQA with Qwen2.5-Math-1.5B, OptPO+PPO achieves better accuracy in all metrics (mean@16, pass@16, and pass@1) than TTRL+PPO while reducing token usage from 262M to 144M, a 45% compute saving. Similar trends are observed for other datasets and RL algorithms, confirming that our early-stopping scheme avoids redundant rollouts on easier samples and filters out harmful samples that occur due to no early-stopping. This ensures that efficiency is significantly improved without harming policy updates.

**Cross-model and cross-algorithm generalization.** Consistent improvements on both Qwen2.5 and Llama backbones indicate that our framework is architecture-agnostic. Even the smaller Llama-3.2-1B shows notable accuracy gains from OptPO-SFT and OptPO over baselines, highlight-

ing that the proposed adaptation mechanism scales robustly across model sizes. Also, our OptPO framework integrates seamlessly with different types of RL algorithms (e.g., PPO, GRPO, Reinforce++) and consistently outperforms TTRL with higher accuracy and substantial efficiency gains, indicating its effectiveness across different RL algorithms.

## 4.3. Evaluation on Test-time Fine-tuning

Table 2 demonstrates our results for TTSFT and OptPO-SFT. Similar to RL-based variants, OptPO-SFT also achieves 15% $\sim$ 60% token consumption savings while maintaining similar or higher accuracy performance metrics than TTSFT.

**SFT-based vs. RL-based variants.** The RL-based OptPO variants generally yield stronger overall accuracy than the OptPO-SFT, particularly on math reasoning benchmarks. This suggests that policy-gradient adaptation better leverages stochastic feedback from rollouts compared to deterministic pseudo-label fine-tuning. Nevertheless, OptPO-SFT achieves more stable improvements over TTSFT across all settings, validating that our adaptive pseudo-labeling significantly enhances SFT-style test-time training.

**Takeaways.** Our results demonstrate that (1) adaptive early-stopping through BSPRT effectively improves efficiency without sacrificing accuracy; (2) OptPO offers the strongest balance between accuracy and compute; and (3) OptPO-SFT provides a lightweight alternative for rapid adaptation when RL-style updates are infeasible. Together, they establish a unified and efficient test-time training paradigm for verifiable reasoning tasks.

## 4.4. Ablation Studies

**Ablation on the number of rollouts used for policy update $N$.** Table 3 shows the performance of TTRL/OptPO when $N$ is 16. OptPO still demonstrates a consistent accuracy performance advantage over the TTRL baseline while saving 43.9% of tokens than the TTRL baseline, which is even higher than 30.0% saving when $N = 32$. This shows the stability of the algorithm, which can still maintain great accuracy with amazing efficiency when we push the rollout budget to its limits.

*Table 3.* Ablation studies on OptPO varying by $N$ for policy update using the Qwen2.5-Math-1.5B backbone.

| Benchmark | AIME 2024 | | | | | AMC | | | | | MATH-500 | | | | | GPQA | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metrics | Mean @16 | Pass @16 | Pass @1 | Toks. | Toks. Saving | Mean @16 | Pass @16 | Pass @1 | Toks. | Toks. Saving | Mean @16 | Pass @16 | Pass @1 | Toks. | Toks. Saving | Mean @16 | Pass @16 | Pass @1 | Toks. | Toks. Saving |
| TTRL-GRPO ($N$=16) | 14.8 | 46.7 | 23.3 | 145M | 20.38% | 46.5 | 83.1 | 54.2 | 162M | 38.40% | 71.4 | 90.8 | 73.6 | 224M | 50.66% | 26.1 | 86.3 | 31.5 | 235M | 66.13% |
| OptPO-GRPO ($N$=16) | 16.5 | 46.7 | 23.3 | 115M | | 46.2 | 83.1 | 55.4 | 100M | | 70.7 | 91.2 | 72.2 | 110M | | 27.1 | 88.3 | 31.5 | 79M | |
| TTRL-GRPO ($N$=32) | 15.6 | 43.3 | 20.0 | 142M | 17.25% | 47.7 | 83.1 | 54.2 | 159M | 26.93% | 72.3 | 90.6 | 73.6 | 218M | 32.57% | 26.5 | 87.8 | 30.5 | 251M | 43.20% |
| OptPO-GRPO ($N$=32) | 14.6 | 46.7 | 20.0 | 117M | | 48.5 | 81.9 | 54.2 | 116M | | 72.3 | 90.8 | 73.8 | 147M | | 26.8 | 87.8 | 34.5 | 143M | |

*Table 4.* Ablation studies on OptPO with different type-I ($\alpha$) and type-II ($\beta$) error budgets using Qwen2.5-Math-1.5B backbone on the MATH-500 benchmark.

| $\alpha/\beta$ | mean@16 | pass@16 | pass@1 | toks. | toks. saving |
|---|---|---|---|---|---|
| 0.01 | 72.0 | 90.8 | 74.8 | 146M | 32.90% |
| 0.03 | 72.1 | 90.4 | 74.2 | 147M | 32.69% |
| 0.05 | 72.3 | 90.8 | 73.8 | 147M | 32.57% |
| 0.07 | 71.9 | 91.0 | 73.8 | 141M | 35.26% |
| 0.1 | 72.2 | 90.6 | 73.8 | 141M | 35.25% |

**Ablations on type-I & type-II error budgets $\alpha$ and $\beta$.** Table 4 shows the performance at OptPO under different type-I and type-II error budgets, when running using Qwen2.5-Math-1.5B on the MATH-500 benchmark. We find that OptPO's performance remains consistently strong and amazingly stable across all configurations. Notably, token savings will improve when we relax error tolerance to 0.07 and 0.1 without accuracy loss. This again confirms the stability and excellent performance of our OptPO framework.

## 5. Conclusion

This paper introduced OptPO, the first statistically grounded framework for efficient test-time policy optimization, which addresses the long-standing inefficiency issue of majority-vote–based rollouts in test-time reinforcement learning for LLMs. By formulating early-stopping consensus estimation as a Bayesian sequential probability ratio test, OptPO adaptively allocates rollouts based on posterior evidence rather than fixed budgets, ensuring reliable pseudo-labels while avoiding redundant sampling. The resulting framework is simple, plug-and-play, and compatible with standard policy optimization methods such as PPO, GRPO, and Reinforce++, as well as supervised test-time fine-tuning. Across diverse benchmark and backbone settings, OptPO consistently preserves or improves accuracy while achieving substantial compute savings. More broadly, this work underscores that efficient, label-free adaptation is both feasible and increasingly essential as LLMs are deployed in unpredictable real-world environments. We hope OptPO serves as a foundation for future research on principled, compute-efficient test-time learning.

## References

Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

Akyürek, E., Damani, M., Zweiger, A., Qiu, L., Guo, H., Pari, J., Kim, Y., and Andreas, J. The surprising effectiveness of test-time training for few-shot learning. *arXiv preprint arXiv:2411.07279*, 2024.

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.

Hardt, M. and Sun, Y. Test-time training on nearest neighbors for large language models. *arXiv preprint arXiv:2305.18466*, 2023.

Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the math dataset. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.

Hu, J. Reinforce++: A simple and efficient approach for aligning large language models. *arXiv preprint arXiv:2501.03262*, 2025.

Hu, J., Zhang, Z., Chen, G., Wen, X., Shuai, C., Luo, W., Xiao, B., Li, Y., and Tan, M. Test-time learning for large language models. *arXiv preprint arXiv:2505.20633*, 2025.

Hübotter, J., Bongni, S., Hakimi, I., and Krause, A. Efficiently learning at test-time: Active fine-tuning of llms. *arXiv preprint arXiv:2410.08020*, 2024.

Li, Y., Lyu, M., and Wang, L. Learning to reason from feedback at test-time. *arXiv preprint arXiv:2502.15771*, 2025.

MAA. American invitational mathematics examination - aime. https://huggingface.co/datasets/AI-MO/aimo-validation-aime, February 2024.

Niu, S., Wu, J., Zhang, Y., Chen, Y., Zheng, S., Zhao, P., and Tan, M. Efficient test-time model adaptation without forgetting. In *International conference on machine learning*, pp. 16888–16905. PMLR, 2022.

Rein, D., Hou, B. L., Stickland, A. C., Petty, J., Pang, R. Y., Dirani, J., Michael, J., and Bowman, S. R. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.

Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y., Wu, Y., et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

Wang, D., Shelhamer, E., Liu, S., Olshausen, B., and Darrell, T. Tent: Fully test-time adaptation by entropy minimization. *arXiv preprint arXiv:2006.10726*, 2020.

Wang, Y., Yang, Q., Zeng, Z., Ren, L., Liu, L., Peng, B., Cheng, H., He, X., Wang, K., Gao, J., et al. Reinforcement learning for reasoning in large language models with one training example. *arXiv preprint arXiv:2504.20571*, 2025.

Zuo, Y., Zhang, K., Sheng, L., Qu, S., Cui, G., Zhu, X., Li, H., Zhang, Y., Long, X., Hua, E., et al. Ttrl: Test-time reinforcement learning. *arXiv preprint arXiv:2504.16084*, 2025.