PLAN THEN ACTION: HIGH-LEVEL PLANNING GUIDANCE REINFORCEMENT LEARNING FOR LLM REASONING

Anonymous authors

000

001

002

004

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032 033 034

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Large language models (LLMs) have demonstrated remarkable reasoning abilities in complex tasks, often relying on Chain-of-Thought (CoT) reasoning. However, due to their autoregressive token-level generation, the reasoning process is largely constrained to local decision-making and lacks global planning. This limitation frequently results in redundant, incoherent, or inaccurate reasoning, which significantly degrades overall performance. Existing approaches, such as tree-based algorithms and reinforcement learning (RL), attempt to address this issue but suffer from high computational costs and often fail to produce optimal reasoning trajectories. To tackle this challenge, we propose Plan-Then-Action Enhanced Reasoning with Group Relative Policy Optimization (*PTA-GRPO*), a two-stage framework designed to improve both high-level planning and fine-grained CoT reasoning. In the first stage, we leverage advanced LLMs to distill CoT into compact high-level guidance, which is then used for supervised fine-tuning (SFT). In the second stage, we introduce a guidance-aware RL method that jointly optimizes the final output and the quality of high-level guidance, thereby enhancing reasoning effectiveness. We conduct extensive experiments on multiple mathematical reasoning benchmarks, including MATH, AIME2024, AIME2025, and AMC, across diverse base models such as Qwen2.5-7B-Instruct, Qwen3-8B, Qwen3-14B, and LLaMA3.2-3B. Experimental results demonstrate that *PTA-GRPO* consistently achieves stable and significant improvements across different models and tasks, validating its effectiveness and generalization.

1 Introduction

Large Language Models (LLMs) have recently demonstrated remarkable reasoning abilities across a wide range of complex tasks (Xu et al., 2025; Plaat et al., 2024; Ke et al., 2025), such as mathematics (Zhang et al., 2024; Wu et al., 2024a; Liu et al., 2023) and programming (Jiang et al., 2024), by leveraging Chain-of-Thought (CoT) reasoning (Wei et al., 2022). Models with strong reasoning capabilities, including Qwen-3 (Yang et al., 2025), DeepSeek-R1 (Wu et al., 2024b), Seed-1.5 thinking (Seed et al., 2025) and GPT-5 thinking (OpenAI, 2025), adopt CoT as a central mechanism to structure their reasoning processes. However, CoT decoding in LLMs is still a token-level Markov Decision Process (MDP) (Ouyang et al., 2022; Wan et al., 2025; Liu et al., 2025): the output of each token is determined by the context sequence generated previously. Under this setting, mainstream decoding is both autoregressive (each decision conditions only on the prefix) and locally greedy (it optimizes short-horizon token likelihood, e.g., via greedy/low-temperature choices). This combination preserves local consistency but offers little global planning, often yielding redundant or drifting chains of thought and propagating early mistakes across long horizons (Yao et al., 2023; Qu et al., 2025; Wan et al., 2025).

Prior work augments LLM reasoning with tree-style algorithms (Zhang et al., 2024; Yao et al., 2023; Wang et al., 2024) such as Monte Carlo Tree Search (Zhang et al., 2024) or heuristic generation tree (Li et al., 2025) to widen exploration beyond single-path decoding. While effective in some cases, these approaches hinge on repeated external queries to the LLM, incurring substantial time and compute (Wang et al., 2024). Crucially, they do not strengthen the model's internal reasoning: performance stems from outside search. When the model cannot verify intermediate steps, the search

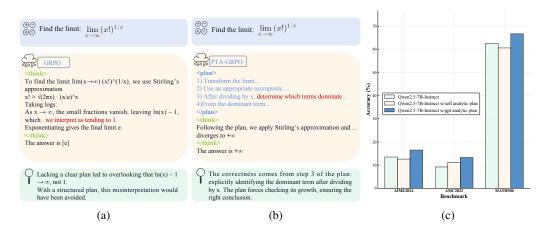


Figure 1: (a) GRPO reasoning processing. (b) *PTA-GRPO* reasoning process. (c) Impact of analytic plan. In (c), the accuracy of different reasoning modes, where Qwen2.5-7B-Instruct is considered as the base model. Yellow indicates the base model using CoT reasoning, blue indicates the base model reasoning with its own self-generated analytic plan, and green indicates the base model reasoning with an analytic plan generated by GPT-o1. More test cases of *PTA-GRPO* are shown in Appendix B.2.

simply amplifies bad branches and collapses (Feng et al., 2023). In parallel, recent works inject reflection or backtracking behaviors via RL (Wan et al., 2025; Wang et al., 2025; Gandhi et al., 2025). Such behaviors can, in principle, re-route trajectories and escape local optima (Gandhi et al., 2025). Yet when triggered on corrupted partial solutions, the model tends to reflect on its own errors, reinforcing them and drifting farther from the correct path. This occurs largely due to the absence of a global plan to guide self-reflection, leaving the model without a reliable mechanism to recover. These limitations motivate a new paradigm that improves internal planning rather than relying on external search or post-hoc self-correction.

Motivated by the way humans tackle complex problems (Kahneman, 2011), where first sketches are made and then executed, it is natural to consider whether LLM reasoning could benefit from a similar paradigm. Specifically, an LLM may first produce a compact and general analytic plan before generating a detailed CoT. Such a plan can provide concise and general global guidance (e.g., subgoal decomposition and task scheduling), and conditioning the CoT on this plan helps mitigate local myopia and reduce redundancy. However, certain weaker LLMs (e.g., Qwen-2.5-7B-Instruct (Bai et al., 2023)) lack the ability to generate high-quality analytic plans. As shown in Fig 1c, the analytic plans generated by Qwen-2.5-7B-Instruct are of insufficient quality, which actually degrades the performance of the resulting CoT and answers, whereas plans generated by the stronger model GPT-o1 lead to significant improvements. These phenomena naturally suggest that a promising direction is to enhance the analytic planning ability of LLMs, as generating high-quality analytic plans can substantially improve their reasoning performance.

To cultivate strong analytic plans, a recent advanced strategy is to exploit the advantages of Reinforcement Learning (RL), e.g., trajectory-level, non-differentiable optimization, enhancing plan quality and alignment with downstream CoT, to achieve reliable, globally guided reasoning. However, under above reasoning paradigm for analytic plan, outcome-based RL with Verifiable Rewards (RLVR) strategies (Shao et al., 2024; Yu et al., 2025; Cui et al., 2025), such as Group Relative Policy Optimization (GRPO) (Shao et al., 2024) or Decoupled Clip and Dynamic Sampling Policy Optimization (DAPO) (Yu et al., 2025), are not entirely suitable. This is because such approaches optimize only for the correctness of the final output while overlooking the quality of the analytic planning and intermediate CoT reasoning as the upper part of Fig 2. Consequently, even poorly planned and executed CoT may receive the same reward as well-structured ones, as long as both yield the correct answer. Such limitations underscore the necessity of developing new RL frameworks that can jointly optimize both the analytic planning and the detailed CoT reasoning processes.

Based on the above analysis, we propose *PTA-GRPO* (*plan-then-action enhanced reasoning with Group Relative Policy Optimization*), a novel two-stage plan-reasoning training framework designed to promote explicit higher-order planning and reasoning abilities. In the first stage, we propose a Planning-Structured Reasoning cold-start approach and leverage an advanced LLM to distill the ground-truth CoT into concise high-level guidance. Recent empirical studies (Gandhi et al., 2025;

Yue et al., 2025; Li et al., 2025) have shown that the reasoning capabilities of pre-trained models are largely established during the initial pre-training phase, which implies that reasoning models are inherently constrained by their base models. These base models lack explicit or autonomous high-quality global planning ability. Therefore, it is necessary to cold-start and cultivate such an initial capability. To this end, the advanced LLM summarizes the CoT by extracting core concepts and generating a refined overview of the reasoning path and conclusions. This high-level guidance thinking, together with the CoT, forms a dataset for high-level guidance-based supervised fine-tuning (SFT), thereby providing a cold-start initialization for subsequent reinforcement learning. In the second stage, we propose a plan reason-guidance aware RL method based on the GRPO algorithm, which has shown strong capabilities in LLM reasoning. Unlike traditional GRPO, which rewards the model based solely on the final response, our method incorporates a sophisticated reward mechanism that evaluates the quality of the high-level guidance thinking generated during the reasoning process. This reward system not only encourages the model to generate accurate final responses but also strengthens its ability to produce effective and precise high-level guidance, thereby enhancing the model's whole reasoning ability. Our main contributions are summarized as follows:

- A novel two-stage plan-reasoning framework: We propose *PTA-GRPO*, a two-stage training framework, including high-level guidance planning and guidance-aware reinforcement learning, to foster explicit higher-order planning and reasoning abilities in LLMs.
- **High-level guidance as supervision signal:** In the supervised fine-tuning stage, we leverage an advanced LLM to transform raw chain-of-thought (CoT) into concise high-level guidance, which is combined with the original CoT, providing stronger initialization for reasoning.
- Plan guidance-aware GRPO with refined reward design: In the reinforcement learning stage, we extend GRPO with a reward mechanism that evaluates not only the correctness of the final response but also the quality of high-level guidance, significantly enhancing overall reasoning effectiveness and robustness.

2 Preliminaries and related work

2.1 REASONING IN LARGE LANGUAGE MODELS

The reasoning of an LLM can be formalized as a token-level Markov Decision Process (MDP) (Ouyang et al., 2022; Wan et al., 2025; Liu et al., 2025), where the state is the context sequence, the action is the next token, and the policy is the model's conditional distribution. Given a question q, a response $\mathfrak{o} = [\mathfrak{o}^1, \dots, \mathfrak{o}^T]$ is sampled step by step from $\pi_{\theta}(\cdot \mid q, \mathfrak{o}^{< t})$. Current inference typically relies on CoT, producing a reasoning chain c and final answer, but this purely autoregressive process lacks global planning, often leading to redundancy and incoherence (Wan et al., 2025).

2.2 GROUP RELATIVE POLICY OPTIMIZATION AND ITS EXTENSIONS

GRPO (Shao et al., 2024), proposed by DeepSeek, enhances LLM reasoning without value models by sampling multiple responses per prompt and using the group average reward as a baseline. This simple mechanism has proven effective in mathematical reasoning, code generation, and QA. Subsequent variants refine GRPO from different perspectives: SRPO (Zhang et al., 2025b) reuses samples via history resampling; DAPO (Yu et al., 2025) filters extreme cases with dynamic sampling; Dr.GRPO (Liu et al., 2025) mitigates length bias; EMPO (Zhang et al., 2025a) optimizes semantic entropy directly; and SEED-GRPO (Seed et al., 2025) integrates entropy as an uncertainty measure for more conservative updates. While these methods substantially improve mathematical reasoning, they do not explicitly target higher-order reasoning abilities.

2.3 MOTIVATION

To address the lack of global guidance in LLM reasoning, which often leads to redundancy or off-topic reasoning, inspired by human thinking habits for complex tasks or problems (Kahneman, 2011; Kahneman & Tversky, 2013), we introduce a concise high-level plan t as an outline before generating the detailed CoT c and its corresponding answer. Formally, the model's output can be represented as $\mathfrak{o}=t,c$, where t provides the overall problem-solving direction without involving concrete computational steps, and c is then generated conditioned on both the question q and the plan t, i.e., $c=\pi_{\theta}(\cdot \mid q,t)$. The CoT c and its final answer are guided by the high-level plan t. This plan-then-reason mechanism equips the reasoning process with global guidance, leading to more concise, and accurate CoT generation.

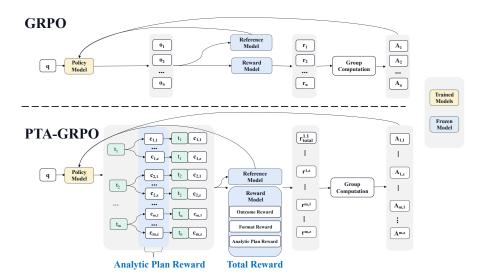


Figure 2: Comparison between GRPO and *PTA-GRPO*. It is worth noting that, to ensure a fair comparison, the number of rollout responses is kept the same between GRPO and PTA-GRPO.

Therefore, in GRPO optimization (the formulas are shown in Appendix B.3) in our study, the objective goes beyond simply ensuring the correctness of the answer in $\mathfrak o$. It also includes enhancing the quality of the high-level plan t, with the aim of producing t more accurately and effectively. By improving t, the model receives structured guidance that can better direct the generation of the CoT c and, consequently, the final answer. This dual focus ensures that the optimization process not only rewards correct answers but also reinforces the production of high-quality intermediate reasoning, leading to more robust and generalizable reasoning behavior.

3 APPROACH OF PTA-GRPO

In this section, we introduce the *PTA-GRPO* training framework, which consists of two key components. (1) Plain Structured Reasoning Cold-Start (PSR-CS). This module serves as a cold-start approach built upon supervised fine-tuning (SFT). Unlike conventional SFT datasets that contain only direct CoT and answers, we first construct a novel dataset that introduces a *general analytical plan* before detailed reasoning. This additional analytical plan provides higher-level guidance, enabling the model to abstract complex problem-solving strategies into concise forms and offering explicit guidance for answer generation. (2) Planning Structure-Guided Reinforcement Learning (PSG-RL). In this stage, we propose a GRPO-based Structure-Guided reinforcement learning algorithm to further enhance the structural reasoning capability of the model. The model is guided to generate general analytical content, whose quality is evaluated and converted into a reward function to determine whether it facilitates more accurate answer generation. This reward signal is then integrated into the GRPO reinforcement learning loop as an explicit optimization objective, thereby forming a closed cycle that continuously improves the effectiveness of the model's reasoning.

3.1 PLANNING STRUCTURED REASONING COLD-START (PSR-CS)

Analytical-Guided SFT Dataset Construction. For LLMs, the ability to perform reasonable planning directly affects whether they can successfully solve a problem. However, existing SFT datasets typically focus only on detailed CoT reasoning and final answers, while neglecting the importance of conducting an overall analytical plan before solving the problem. To address this gap, we propose an analytical-guided dataset, which consists of three components: the problem, a general analytical plan, and the corresponding detailed CoT reasoning with the final answer. This dataset not only injects concise and effective general analytical knowledge into LLMs to provide an overall problem-solving perspective but also trains them to transform such general plans into concrete reasoning processes, thereby enhancing their overall reasoning capabilities. Formally, we define the dataset as $D_{\text{PSR-CS}} = \{q_i, t_i, c_i\}_{i=1}^n$, which contains n tuples, where each tuple comprises the problem q_i , the general analytical plan t_i , and the corresponding detailed reasoning with the final answer c_i . In our constructed dataset, the general analytical plan t_i is enclosed within the plan > ... < plan> tags, which clearly distinguishes the high-level problem-solving idea.

Meanwhile, the specific response c_i is further structured: the chain-of-thought (CoT) is wrapped in <think>...

think>...
, and the final answer is wrapped in <answer>...</answer>, thereby providing a hierarchical representation of planning, reasoning, and answering.

In practice, we sampled 10K instances from the Openthoughts (Guha et al., 2025) dataset as our base. Openthoughts is a large-scale open reasoning dataset that covers a wide range of problems along with their detailed CoT reasoning processes. We then employed the powerful open-source reasoning model Qwen3-235B (Yang et al., 2025) as the teacher model. For each instance, we input the problem q_i and its detailed reasoning c_i into the advanced model to generate the corresponding general analytical plan t_i . Through this process, we distilled general analytical knowledge from a strong LLM and injected it into our target models to enhance their overall reasoning capability.

SFT-based Cold-Start Initialization Optimization. At this stage, we aim to inject structured reasoning capabilities into the initial policy model π_{θ} through SFT, which serves as an effective way to expand the knowledge and abilities of LLMs (Shah et al., 2025). Specifically, we optimize the model parameters by minimizing the discrepancy between the model outputs and the reference outputs provided in the analytical-guided dataset D_{SRCS} , thereby enabling the model to gradually acquire structured reasoning patterns. The fine-tuning process can be formulated as:

$$\theta_{\text{SFT}} = \min_{\theta} \quad \mathbb{E}_{(q_i, t_i, c_i) \in \mathcal{D}_{\text{SRCS}}} \left[\sum_{i=1}^{n} \log \left(\pi_{\theta}(t_i, c_i \mid q_i) \right) \right]. \tag{1}$$

 $\theta_{\rm SFT}$ refers to the parameter set learned through supervised fine-tuning on the analytical-guided dataset. Based on these optimized parameters, $\pi_{\theta_{\rm SFT}}$ denotes the resulting policy model that embodies structured reasoning capabilities. By explicitly injecting high-level analytical plans before detailed CoT reasoning, the policy model is guided to generate solutions in a more systematic and interpretable manner.

3.2 PLAN STRUCTURE-GUIDED REINFORCEMENT LEARNING (PSG-RL)

After obtaining the policy model $\pi_{\theta_{\rm SFT}}$ from the SFT stage, the RL phase then focuses on improving the model's planning capability and ensuring its effective execution. At this stage, we not only consider the correctness of CoT c and its answer as part of the reward signal, but also evaluate the quality of the analytical plan t, which is incorporated as another important aspect of the reward signal.

3.2.1 ANALYTICAL PLAN-GUIDED REWARD AUGMENTATION IN GRPO

In PTA-GRPO, we design a composite reward function that integrates three aspects: the analytical planning reward ($r_{\rm analytical}$) to encourage structured reasoning plans, the outcome accuracy reward ($r_{\rm outcome}$) to ensure correct final results, and the structured format reward ($r_{\rm format}$) to enforce clear and consistent output. Together, these rewards are combined into the total reward $R_{\rm total}$, which enhances the model's planning capability, reasoning accuracy, and response reliability.

Analytical Plan Reward. Since directly evaluating the quality of an analytical plan t is difficult in practice, we instead use computable and optimizable surrogate objectives to measure the probability that it guides a specific CoT reasoning process toward the correct answer, where a higher probability intuitively reflects a higher-quality plan. Based on this insight, we design the reward for the analytical plan $r_{\rm analytic}$, which is defined by the probability that the analytical plan can guide a CoT reasoning process toward the correct answer. To achieve the above goal, we construct a response group G through a two-step process. Given a question q, the policy model first samples a set of m candidate analytical plans $\{t_i\}_{i=1}^m$, where $t_i \sim \pi_{\theta}(\cdot \mid q)$ and each analytical plans t_i is a concise, text-based outline of how to approach q. Then, for each analytical plan t_j , following (Lu et al., 2025), we resample z detailed CoT $\{c_{i,k}\}_{k=1}^z$ under guidance of t_i , where each $c_{i,k}$ is drawn as $c_{i,k} \sim \pi_{\theta}(\cdot \mid t_i, q)$. The response group G consists of m analytical plans, each associated with z CoT, where $G = \left\{\{(t_i, c_{i,k})\}_{k=1}^z\right\}_{i=1}^m$. For each response from G can be regarded as planning-CoT paris, and the reward $r_{\rm analytic}$ assigned to t_i is defined as the empirical accuracy of its resampled outcomes:

$$r_{\text{analytic}}(t_i) = \text{Softmax}\left(\frac{1}{z}\sum_{k=1}^{z} \mathbb{I}[\hat{y}_{i,k} = y]\right),$$
 (2)

where $\mathbb{I}[\cdot]$ is the indicator function, $\hat{y}_{i,k}$ denotes the final expected answer extracted from $c_{i,k}$, and y is the ground-truth answer of q. Through the policy model driven by $r_{\text{analytic}}(\cdot)$, more accurate analytic plans t can be generated, thereby improving the probability of obtaining the correct prediction $\Pr(\hat{y} = y \mid t, q)$. In addition, we apply the Softmax to exponentially amplify the differences between scores, making high-scoring planning more prominent while further suppressing low-scoring ones.

Outcome Reward. The outcome reward, defined as r_{outcome} , is a result-based terminal reward similar to GRPO, used to evaluate whether the predicted answer aligns with the ground truth. For each plan–CoT response $(t_i, c_{i,k})$, the outcome reward r_{outcome} is defined as follows:

$$r_{\text{outcome}} = \begin{cases} 1, & \hat{y}_{i,k} = y, \\ 0, & \text{else.} \end{cases}$$
 (3)

The outcome reward r_{outcome} encourages the policy model to learn to follow the analytical plan t_i and to develop the ability to generate answers that strive for correctness.

Format Reward. The format reward $r_{\rm format}$ is designed to regulate the overall structure of the model response, ensuring both conformity to the desired format and control over the output length. It consists of two components: $r_{\rm structure}$ and $r_{\rm length}$. Specifically, $r_{\rm structure}$ enforces that the policy model's response adheres to the predefined structural template, i.e., <plan>...</plan>, <think>...</think>, and <answer>...</answer>. Meanwhile, $r_{\rm length}$ serves as an auxiliary reward that encourages the model to generate concise and efficient token sequences, thereby reducing redundant or uninformative content.

To provide a clearer illustration of each reward, we present its detailed formulation as follows. We begin with the format reward r_{format} , which is defined as:

$$r_{\rm format} = \begin{cases} 0.2, & \text{if the response strictly follows the predefined template} \\ 0, & \text{otherwise.} \end{cases} \tag{4}$$

This function enforces a binary constraint on the output structure: a full reward is granted only when the response strictly adheres to the predefined template, thereby ensuring the consistency and parsability of the generated results.

For response length, the optimal number of tokens varies across different questions, making it difficult to predefine a fixed target length. Therefore, for all responses generated for a given question, we select the shortest correct response length as the reference length T, defined as:

$$T = \min\{ |\{t_i, c_{i,k}\}| \mid \hat{y}_{i,k} = y \}, \tag{5}$$

where $|\{t_i, c_{i,k}\}|$ denotes the token length of response $\{t_i, c_{i,k}\}$. Here, T represents the shortest executable token length required to obtain the correct answer to a given question. It can be regarded as the optimal reference length under current knowledge, toward which other correct responses should converge in order to minimize redundancy while preserving correctness. For each response $\{t_i, c_{i,k}\} \in G$, the length reward r_{length} can be expressed as:

$$r_{\text{length}}(\{t_i, c_{i,k}\}) = \alpha \cdot \exp(-\frac{||\{t_i, c_{i,k}\}| - T|}{T_{\text{max}} - T}),$$
 (6)

where α is not a hyperparameter, and $T_{\rm max}$ does not denote the maximum output length set for the policy model. The reward becomes larger as the response length approaches the reference length T, encouraging the model to generate concise yet correct responses.

The format reward r_{format} , defined as $r_{\text{format}} = r_{\text{structure}} + r_{\text{length}}$, ensures that the output not only adheres to the required format, but also guarantees the conciseness of the output response.

Total Reward. The above three rewards together constitute the total reward R_{total} for each response as:

$$R_{\text{total}} = R_{\text{analytic}} + \beta \cdot R_{\text{outcome}} + R_{\text{format}}, \tag{7}$$

where β represents the hyperparameter. We first obtain a total reward set $\{\{r_{total}^{i,k}\}_{i=1}^m\}_{k=1}^z$, where $r_{total}^{i,k}$ denotes the total reward of the k-th CoT generated under the guidance of the i-th analytic. Based on this reward, we compute the corresponding advantage function $A_{i,k}$ using Eq. 9, and subsequently incorporate it into the update rule in Eq. 8 to optimize the model.

Advantages Compared with Conventional GRPO. Compared with standard GRPO, which primarily relies on sparse task-level accuracy supervision, our guidance-aware PTA-GRPO framework introduces several critical improvements. First, by incorporating the $r_{\rm analytic}$ indicator, the model strengthens its analytic planning ability, leveraging self-generated plans to guide subsequent computation. Second, the outcome reward $r_{\rm outcome}$ encourages the policy model to follow the analytic plan and enhance its reasoning capability under such structured guidance. Third, the format reward $r_{\rm format}$ promotes stable and standardized reasoning patterns, optimizing outputs toward being both minimal response length and correct. Together, these enhancements enable PTA-GRPO to achieve stronger high-level analytic planning and improved reasoning performance in complex tasks compared to standard GRPO.

3.3 THEORETICAL PERFORMANCE ANALYSIS

In this section, we theoretically analyze the impact of optimizing r_{analytic} on the performance of the policy model on the probability of errors. Our theoretical findings are as follows.

Theorem 3.1. Let q denote the input question, t the analytic plan, \hat{y} the answer predicted by the policy model, and y the ground-truth answer. With error probability p_{error} , it holds that:

$$p_{error} \le \frac{1}{2} [H(y) - I(\hat{y}, y \mid t, q)], \quad p_{error} = \Pr(y \ne \hat{y}),$$

where $H(\cdot)$ denotes the entropy, and $I(\cdot)$ denotes the mutual information.

The proof can be seen in appendix B.4. Leveraging the conclusion from (Qian et al., 2025), since H(y) depends solely on the fixed distribution of the answer and is independent of the model's reasoning steps, it can therefore be regarded as a constant. In our Theorem 3.1, the upper bound of the error probability $p_{\rm error}$ is governed by the conditional mutual information $I(\hat{y};y\mid t,q)$, which measures the statistical dependence between the predicted output \hat{y} and the true label y, given the auxiliary analytic plan t. In other words, the larger the shared information between \hat{y} and y under the guidance of the analytic plan t, the tighter the achievable upper bound on the error probability. Thus, our theoretical analysis illustrates that enabling the LLM to generate an analytic plan t is essential for improving reasoning performance.

Remark 3.2. By the definition of mutual information, $I(\hat{y};y\mid q,t)=H(y\mid q,t)-H(y\mid \hat{y},q,t).$ Note that $H(y\mid q,t)$ is solely determined by the underlying data distribution of (q,t,y) and is independent of the model's prediction \hat{y} . Hence, $H(y\mid q,t)$ can be regarded as a constant with respect to the learning or inference process. Therefore, increasing the mutual information $I(\hat{y};y\mid q,t)$ essentially amounts to reducing $H(y\mid \hat{y},q,t)$, i.e., making the true answer y less uncertain once the model prediction \hat{y} is observed. Theorem 3.1 demonstrates that optimizing $r_{\rm analytic}$ effectively reduces the error probability of the policy model.

4 EXPERIMENT

Based Models. To evaluate *PTA-GRPO*, we adopt four base models of varying scales and series: LLaMA3.2-3B (Dubey et al., 2024), Qwen2.5-7B-Instruct (Bai et al., 2025), Qwen3-8B, and Qwen3-14B (Yang et al., 2025), enabling a comprehensive assessment of its robustness across architectures. Training details are in Section B.5.

Training Datasets and Benchmarks. For SFT, we use 10K samples from Openthoughts (Guha et al., 2025) with injected planning knowledge (Section 3.1). For RL, we sample 14K problems from DeeMath (He et al., 2025), which offers graded difficulty and is rigorously decontaminated to avoid benchmark leakage. We evaluate on AIME24, AIME25, MATH500, AMC23, Minerva, and Olympiad, reporting average accuracy over 16 independent runs.

Baseline. We compare *PTA-GRPO* with the base model, GRPO (Shao et al., 2024), and DAPO (Yu et al., 2025). For fairness, all methods use the same SFT and RL data (differing only in the improved SFT portion), and the RL stage maintains an equal number of responses.

4.1 PERFORMANCE OF PTA-GRPO

Table 1 shows that our method (*PTA-GRPO*) consistently outperforms both the base models and other RLVR approaches across different model scales and evaluation benchmarks. For relatively weaker

Table 1: Performance comparison of different RLVR methods using various base models. **Bold** is best per block.

Method	MATH500	AIME24	AIME25	AMC23	Average
Qwen2.5-7B-Instruct	62.40	12.24 3.52		52.75	32.73
GRPO	82.74	27.52	22.33	63.59	49.04
DAPO	83.92	28.90 21.25		67.75	50.45
PTA-GRPO	85.57	30.26	25.97	70.24	53.01
LLaMA3.2-3B	34.27	3.33	2.74	18.75	14.77
GRPO	55.19	16.27	14.22	38.25	30.98
DAPO	54.27	18.35	16.53	38.25	31.85
PTA-GRPO	60.25	20.50	14.27	40.37	33.85
Qwen3-8B	90.27	66.67	51.53	90.05	74.63
GRPO	92.93	68.27	54.23	91.97	76.85
DAPO	91.27	66.39	50.08	91.33	74.77
PTA-GRPO	93.31	68.88	54.29	92.29	77.19
Qwen3-14B	91.27	72.65	70.03	94.33	82.07
GRPO	90.28	71.29	71.29	94.92	81.95
DAPO	91.07	72.33	70.92	95.20	82.38
PTA-GRPO	91.93	73.90	71.55	94.97	83.09

backbones such as Qwen2.5-7B-Instruct and LLaMA3.2-3B, *PTA-GRPO* delivers the most significant improvements, raising the average scores by over 20 points compared to the raw models and further surpassing GRPO and DAPO by clear margins. The gains are particularly notable on challenging tasks like AIME25, where *PTA-GRPO* achieves improvements of up to four points compared to the strongest baseline.

For stronger backbones such as Qwen3-8B and Qwen3-14B, the room for improvement is smaller since the base performance is already high. Nevertheless, *PTA-GRPO* still provides consistent and measurable gains across nearly all benchmarks, establishing new best results on average scores. Importantly, the method shows no signs of degradation and demonstrates robust generalization across tasks, making it effective not only for improving weaker models but also for further pushing the limits of state-of-the-art models.

4.2 IMPACT OF RL DATA SCALING

Table 2 shows how the performance of Qwen2.5-7B-Instruct on four math benchmarks changes as the RL data scale increases from 4k to 14k. Overall, all tasks steadily improve with larger data sizes, with the average score rising from 48.94 to 53.01, indicating consistent gains from more training data. Specifically, MATH500 remains the strongest across all scales (82.27→85.57), while AIME24 and AIME25, though starting lower, achieve the largest relative improvements, particularly AIME25, which increases from 21.03 to 25.97, a gain of over 23

Table 2: Impact of data scale of RL on *PTA-GRPO*, where Qwen2.5-7B-Instruct is considered as base model. **Bold** is best per block.

Data scale	MATH500	AIME24	AIME25	AMC23	Average
4k	82.27	27.22	21.03	65.22	48.94
8k	83.59	28.23	22.29	68.29	50.60
11k	84.23	29.33	24.51	69.37	51.86
14k	85.57	30.26	25.97	70.24	53.01

4.3 ABLATION ANALYSIS

The results from the two tables demonstrate both the effect of data scale and the importance of different components in *PTA-GRPO*. As shown in Table 2, increasing the RL data size from 4k to 14k steadily improves performance across all benchmarks, with the average score rising from 48.94 to 53.01. Notably, harder tasks such as AIME24 and AIME25 benefit the most, while the trend suggests that performance has not yet saturated at 14k. Table 3 presents the ablation analysis, where removing SFT leads to a significant drop in average performance (41.74), highlighting its necessity. Excluding the format reward slightly boosts AIME24 but lowers the overall average to 52.34, while excluding the analytic reward reduces the average further to 49.86, indicating its critical role in enhancing reasoning quality. The complete *PTA-GRPO* achieves the best overall performance

(53.01), confirming that the combination of SFT, format reward, and analytic reward is essential for maximizing both stability and accuracy.

Table 3: Ablation analysis on *PTA-GRPO*, where Qwen2.5-7B-Instruct is considered as base model. **Bold** is best per block.

Method	MATH500	AIME24	AIME25	AMC23	Average
PTA-GRPO w/o SFT	79.25	16.25	12.25	59.22	41.74
PTA - $GRPO$ _{w/or_{format}}	85.37	31.23	24.52	68.25	52.34
PTA-GRPO W/o ranalytic	81.03	28.22	23.85	66.33	49.86
PTA-GRPO	85.57	30.26	25.97	70.24	53.01

4.4 IMPACT OF ANALYTIC PLAN ON SFT

Table 4: The impact of datasets containing analytic planning on SFT. **Bold** is best per block.

Base Model	Method	MATH500	AIME24	AIME25	AMC23	Average
Qwen2.5-7B-Instruct	SFT w/o planning	78.28	21.66	19.66	60.53	45.03
	SFT w/ planning	80.40	25.25	20.33	63.75	47.43
Qwen3-8B	SFT w/o planning	91.02	70.03	50.25	92.39	75.92
	SFT w/ planning	92.53	71.97	51.77	93.55	77.46

In Table 4, we compare standard SFT (SFT w/o planning) with SFT using \mathcal{D}_{SRCS} augmented by analytic plans (SFT w/ planning). The results demonstrate that incorporating analytic plans consistently yields improvements across all tasks and models. For Qwen2.5-7B-Instruct, the average score increases from 45.03 to 47.43, with gains of 0.67–3.59 points across the four benchmarks, highlighting the stronger reliance of smaller models on external planning signals. For the more capable Qwen3-8B, although it already possesses stronger reasoning ability, analytic plans still improve the average score from 75.92 to 77.46, with gains mainly between 1–2 points. Overall, analytic plans serve as structured reasoning supervision signals that not only significantly enhance the reasoning ability of smaller models but also provide consistent fine-grained improvements for larger models.

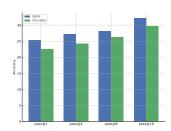


Figure 3: Effect of scaling test-time compute on AIME25 (Pass@K), with Qwen2.5-7B-Instruct as the base model.

4.5 RESULTS OF SCALING TEST-TIME

We next examine the effectiveness of *PTA-GRPO* under multiple sampling at test time. As shown in Fig. 3, *PTA-GRPO* consistently outperforms GRPO on the AIME2025 dataset across Pass@1, Pass@4, Pass@8, and Pass@16. This demonstrates that *PTA-GRPO* maintains high precision under low-sample conditions, while further exhibiting stronger solution coverage as the number of samples increases.

4.6 TRAINING DYNAMICS OF PTA-GRPO

Appendix B.1 Fig. 4 and Fig. 5 illustrate the training dynamics of QWEN3-8B and QWEN2.5-7B-Instruct, respectively. As shown in the figures, our method outperforms GRPO in terms of accuracy reward and response length, indicating the effectiveness of the introduced component. It is worth noting that in Fig. 4 (b), our approach achieves lower entropy compared to GRPO. This suggests that for stronger models, our method encourages the development of more reasonable analytic plans, enabling the model to complete a given trajectory with greater confidence and ultimately achieving higher accuracy.

5 CONCLUSION

We propose Plan-Guide Enhanced Reasoning with Group Relative Policy Optimization (*PTA-GRPO*), which integrates high-level planning with fine-grained reasoning to alleviate the lack of global planning in traditional CoT reasoning. Experimental results show that *PTA-GRPO* achieves stable and significant improvements across multiple mathematical reasoning benchmarks and model scales, validating its effectiveness and generalizability.

6 ETHICS STATEMENT

This research has been conducted in alignment with the ICLR Code of Ethics. We are committed to responsible stewardship of machine learning research, ensuring that our work advances knowledge while considering its potential societal impacts. In particular, we uphold high standards of scientific rigor, transparency, and reproducibility, and we affirm that no data has been falsified, fabricated, or misrepresented. Our study avoids harm by carefully considering possible negative consequences and by respecting privacy, fairness, and inclusiveness in the use of data and methods. All data used complies with relevant ethical approvals and license requirements, and precautions have been taken to prevent re-identification or misuse. We respect the intellectual contributions of others and provide appropriate credit where due. We believe this work contributes positively to human well-being by addressing problems of scientific and social relevance in ways that are transparent, responsible, and consistent with the principles of the ICLR Code of Ethics.

7 REPRODUCIBILITY STATEMENT

We have taken several steps to ensure the reproducibility of our work. The main experimental setup, including model architectures, training procedures, and evaluation metrics, is described in detail in the main paper and appendix. To facilitate reproducibility, we will release the majority of the code with an anonymous code link (shown in the Appendix) during the review process. If the paper is accepted, we commit to releasing the complete code base for all major experiments, along with detailed documentation and instructions for reproducing the reported results.

REFERENCES

- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 1(2):3, 2023.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang, Wendi Li, Bingxiang He, Yuchen Fan, Tianyu Yu, Qixin Xu, Weize Chen, et al. Process reinforcement through implicit rewards. *arXiv preprint arXiv:2502.01456*, 2025.
- 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 e-prints*, pp. arXiv–2407, 2024.
- Xidong Feng, Ziyu Wan, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, and Jun Wang. Alphazero-like tree-search can guide large language model decoding and training. *arXiv* preprint arXiv:2309.17179, 2023.
- Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh, Nathan Lile, and Noah D Goodman. Cognitive behaviors that enable self-improving reasoners, or, four habits of highly effective stars. *arXiv* preprint arXiv:2503.01307, 2025.
- Etash Guha, Ryan Marten, Sedrick Keh, Negin Raoof, Georgios Smyrnis, Hritik Bansal, Marianna Nezhurina, Jean Mercat, Trung Vu, Zayne Sprague, et al. Openthoughts: Data recipes for reasoning models. *arXiv preprint arXiv*:2506.04178, 2025.
- Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu Chen, Yue Wang, Linfeng Song, Dian Yu, Zhenwen Liang, Wenxuan Wang, et al. Deepmath-103k: A large-scale, challenging, decontaminated, and verifiable mathematical dataset for advancing reasoning. *arXiv preprint arXiv:2504.11456*, 2025.
- Juyong Jiang, Fan Wang, Jiasi Shen, Sungju Kim, and Sunghun Kim. A survey on large language models for code generation. *arXiv preprint arXiv:2406.00515*, 2024.
- Daniel Kahneman. Thinking, fast and slow. Farrar, Straus and Giroux, 2011.
- Daniel Kahneman and Amos Tversky. Prospect theory: An analysis of decision under risk. In *Handbook of the fundamentals of financial decision making: Part I*, pp. 99–127. World Scientific, 2013.
- Zixuan Ke, Fangkai Jiao, Yifei Ming, Xuan-Phi Nguyen, Austin Xu, Do Xuan Long, Minzhi Li, Chengwei Qin, Peifeng Wang, Silvio Savarese, et al. A survey of frontiers in llm reasoning: Inference scaling, learning to reason, and agentic systems. *arXiv preprint arXiv:2504.09037*, 2025.
- Yizhi Li, Qingshui Gu, Zhoufutu Wen, Ziniu Li, Tianshun Xing, Shuyue Guo, Tianyu Zheng, Xin Zhou, Xingwei Qu, Wangchunshu Zhou, et al. Treepo: Bridging the gap of policy optimization and efficacy and inference efficiency with heuristic tree-based modeling. *arXiv preprint arXiv:2508.17445*, 2025.
- Yixin Liu, Avi Singh, C Daniel Freeman, John D Co-Reyes, and Peter J Liu. Improving large language model fine-tuning for solving math problems. *arXiv preprint arXiv:2310.10047*, 2023.
- 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.
 - Fanbin Lu, Zhisheng Zhong, Shu Liu, Chi-Wing Fu, and Jiaya Jia. Arpo: End-to-end policy optimization for gui agents with experience replay. *arXiv preprint arXiv:2505.16282*, 2025.

- OpenAI. Gpt-5 system card. Technical report, 2025. URL https://cdn.openai.com/gpt-5-system-card.pdf. Accessed: 2025-08-13.
 - 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.
 - Aske Plaat, Annie Wong, Suzan Verberne, Joost Broekens, Niki van Stein, and Thomas Bäck. Reasoning with large language models, a survey. *CoRR*, 2024.
 - Chen Qian, Dongrui Liu, Haochen Wen, Zhen Bai, Yong Liu, and Jing Shao. Demystifying reasoning dynamics with mutual information: Thinking tokens are information peaks in llm reasoning. *arXiv* preprint arXiv:2506.02867, 2025.
 - Xiaoye Qu, Yafu Li, Zhaochen Su, Weigao Sun, Jianhao Yan, Dongrui Liu, Ganqu Cui, Daizong Liu, Shuxian Liang, Junxian He, et al. A survey of efficient reasoning for large reasoning models: Language, multimodality, and beyond. *arXiv preprint arXiv:2503.21614*, 2025.
 - John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
 - ByteDance Seed, Jiaze Chen, Tiantian Fan, Xin Liu, Lingjun Liu, Zhiqi Lin, Mingxuan Wang, Chengyi Wang, Xiangpeng Wei, Wenyuan Xu, et al. Seed1. 5-thinking: Advancing superb reasoning models with reinforcement learning. *arXiv preprint arXiv:2504.13914*, 2025.
 - Darsh J Shah, Peter Rushton, Somanshu Singla, Mohit Parmar, Kurt Smith, Yash Vanjani, Ashish Vaswani, Adarsh Chaluvaraju, Andrew Hojel, Andrew Ma, et al. Rethinking reflection in pretraining. *arXiv preprint arXiv:2504.04022*, 2025.
 - 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.
 - Zhongwei Wan, Zhihao Dou, Che Liu, Yu Zhang, Dongfei Cui, Qinjian Zhao, Hui Shen, Jing Xiong, Yi Xin, Yifan Jiang, et al. Srpo: Enhancing multimodal llm reasoning via reflection-aware reinforcement learning. *arXiv* preprint arXiv:2506.01713, 2025.
 - Ante Wang, Linfeng Song, Ye Tian, Baolin Peng, Dian Yu, Haitao Mi, Jinsong Su, and Dong Yu. Litesearch: Efficacious tree search for llm. *arXiv preprint arXiv:2407.00320*, 2024.
 - Haozhe Wang, Chao Qu, Zuming Huang, Wei Chu, Fangzhen Lin, and Wenhu Chen. Vl-rethinker: Incentivizing self-reflection of vision-language models with reinforcement learning. *arXiv* preprint *arXiv*:2504.08837, 2025.
 - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837, 2022.
 - Yiran Wu, Feiran Jia, Shaokun Zhang, Hangyu Li, Erkang Zhu, Yue Wang, Yin Tat Lee, Richard Peng, Qingyun Wu, and Chi Wang. Mathchat: Converse to tackle challenging math problems with llm agents. *ICLR 2024 Workshop on LLM Agents*, 2024a.
 - Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, et al. Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding. *arXiv preprint arXiv:2412.10302*, 2024b.
 - Fengli Xu, Qianyue Hao, Zefang Zong, Jingwei Wang, Yunke Zhang, Jingyi Wang, Xiaochong Lan, Jiahui Gong, Tianjian Ouyang, Fanjin Meng, et al. Towards large reasoning models: A survey of reinforced reasoning with large language models. *arXiv preprint arXiv:2501.09686*, 2025.
 - An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.

- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.
- Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. arXiv preprint arXiv:2503.14476, 2025.
- Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Shiji Song, and Gao Huang. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model? *arXiv* preprint arXiv:2504.13837, 2025.
- Di Zhang, Jianbo Wu, Jingdi Lei, Tong Che, Jiatong Li, Tong Xie, Xiaoshui Huang, Shufei Zhang, Marco Pavone, Yuqiang Li, et al. Llama-berry: Pairwise optimization for o1-like olympiad-level mathematical reasoning. *arXiv preprint arXiv:2410.02884*, 2024.
- Qingyang Zhang, Haitao Wu, Changqing Zhang, Peilin Zhao, and Yatao Bian. Right question is already half the answer: Fully unsupervised llm reasoning incentivization. *arXiv preprint arXiv:2504.05812*, 2025a.
- Xiaojiang Zhang, Jinghui Wang, Zifei Cheng, Wenhao Zhuang, Zheng Lin, Minglei Zhang, Shaojie Wang, Yinghan Cui, Chao Wang, Junyi Peng, et al. Srpo: A cross-domain implementation of large-scale reinforcement learning on llm. *arXiv preprint arXiv:2504.14286*, 2025b.

A THE USE OF LARGE LANGUAGE MODELS

In preparing this manuscript, we used a Large Language Model (LLM) solely to assist with minor language polishing and improvements in readability. The LLM did not contribute to research ideation, analysis, or substantive writing. All scientific content and conclusions are entirely the responsibility of the authors.

B APPENDIX

B.1 TRAINING DYNAMICS

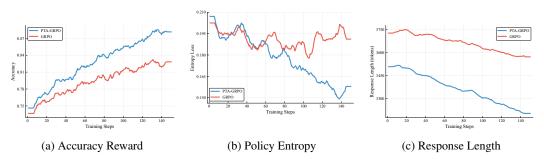


Figure 4: Training Dynamics of PTA-GRPO with Qwen3-8B.

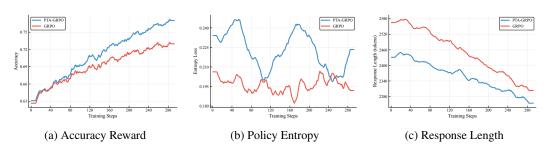


Figure 5: Training Dynamics of *PTA-GRPO* with Qwen2.5-7B-Instruct.

B.2 TEST CASE OF PTA-GRPO

Generated Samples in Real Test cases of PTA-GRPO (Sample 1)

Question: Determine the interval of convergence for the series

$$1 + \frac{x}{2} + \frac{2!x^2}{3^2} + \frac{3!x^3}{4^3} + \dots$$

<plan>

1. Identify the general term of the series. 2. Apply the Ratio Test to find the radius of convergence. 3. Check the endpoints of the interval separately to determine convergence or divergence there. 4. State the final interval of convergence.

 We start by writing down the general term of the series. The pattern is clear: the n-th term looks like

$$\frac{n!}{(n+1)^n}x^n.$$

To test for convergence, the Ratio Test is natural. Taking the ratio a_{n+1}/a_n , most factors cancel, and it reduces to

$$|x|\left(\frac{n+1}{n+2}\right)^{n+1}.$$

As $n \to \infty$, the power term tends to e^{-1} . So the limit is |x|/e. By the Ratio Test, convergence requires |x|/e < 1, which means |x| < e. That gives us the radius of convergence R = e.

Now check the endpoints. At x = e, the terms are $\frac{n!}{(n+1)^n}e^n$. Stirling's formula shows that $n! \sim (n/e)^n \sqrt{2\pi n}$, so the terms behave like \sqrt{n} times a constant. They don't go to zero, so the series diverges. At x=-e, the alternating sign doesn't help because the size of the terms again grows like \sqrt{n} . They still don't go to zero, so the series diverges there as well. Therefore, the series converges only strictly inside the interval, and the interval of convergence is

$$(-e,e)$$
.

Answer:

$$(-e, e)$$

Generated Samples in Real Test cases of PTA-GRPO (Sample 2)

Question: What is the least possible value of the sum |x-1| + |x-1.5| + |x-2|?

1. Recognize that this is a minimization problem involving sums of absolute values. 2. Recall that for a sum of absolute values of linear functions, the minimum occurs at a median of the points. 3. Identify the points: 1, 1.5, and 2. 4. Determine the median of these three points. 5. Calculate the sum at the median point to find the minimum value.

We want to minimize

$$f(x) = |x - 1| + |x - 1.5| + |x - 2|.$$

General fact: For $|x - a_1| + \cdots + |x - a_n|$, the minimum occurs when x is the median of

Here, the three points are 1, 1.5, 2. Ordered: 1 < 1.5 < 2. So the median is 1.5. Evaluate at x = 1.5:

$$|1.5 - 1| + |1.5 - 1.5| + |1.5 - 2| = 0.5 + 0 + 0.5 = 1.$$

Check at neighbors:

* At x = 1: 0 + 0.5 + 1 = 1.5. * At x = 2: 1 + 0.5 + 0 = 1.5. * At x = 1.25: 0.25 + 0.25 + 0.75 = 1.25. * At x = 1.75: 0.75 + 0.25 + 0.25 = 1.25.

All larger than 1. Thus, the minimum occurs uniquely at x = 1.5, with value 1.

Answer:

Generated Samples in Real Test cases of *PTA-GRPO* (Sample 3)

Question: Find the range of p for which the stability condition $|\xi| \le 1$ holds, given that

$$|\xi|^2 = 1 - 4p^2(1 - p^2)s^4$$
, $s = \sin\left(\frac{\omega}{2}\right) \in [0, 1]$.

1. Understand the stability condition: $|\xi| \le 1$. 2. Substitute the given expression for $|\xi|^2$. 3. The stability condition $|\xi| \le 1$ is equivalent to $|\xi|^2 \le 1$. 4. Since $|\xi|^2$ is expressed as $1-4p^2(1-p^2)s^4$, ensure it is non-negative (≥ 0) to avoid complex values. 5. Analyze the sign of $p^2(1-p^2)$. 6. Determine the condition on p such that $|\xi|^2 \in [0,1]$ for all $s \in [0,1]$. 7. Conclude the range of p.

We have

$$|\xi|^2 = 1 - 4p^2(1 - p^2)s^4, \quad s \in [0, 1].$$

For stability, we need $|\xi| \le 1$, i.e., $|\xi|^2 \le 1$. That means

$$1 - 4p^2(1 - p^2)s^4 \le 1.$$

This simplifies to

$$-4p^2(1-p^2)s^4 \le 0 \implies p^2(1-p^2) \ge 0.$$

Since $p^2 \ge 0$, this requires $1-p^2 \ge 0 \implies p^2 \le 1$. Thus $|p| \le 1$. Additionally, we must ensure $|\xi|^2 \ge 0$. For $|p| \le 1$, the term $4p^2(1-p^2)s^4 \ge 0$, so $|\xi|^2 = 1$ – (nonnegative) $\in [0,1]$. This guarantees both realness and the upper bound condition.

If |p| > 1, then $1 - p^2 < 0$, so $p^2(1 - p^2) < 0$. In this case the subtracted term becomes negative, effectively making $|\xi|^2 > 1$ for some s, which breaks stability.

Hence the stability condition holds exactly when

$$p \in [-1, 1].$$

Answer:

-1, 1

PTA-GRPO prompt

{{ content | trim }}

Begin by outlining a high-level plan wrapped in <plan></plan> tags.

- This plan should capture only the major phases, strategic choices, and conditional branches.
- Avoid low-level steps, calculations, or detailed reasoning here. Next, reason step by step within <think></think>.
- During reasoning, critically evaluate the initial plan. you find any errors, inconsistencies, or improvements needed, revise your plan mentally and continue reasoning based on the revised plan.

- Explicitly state if you are revising the plan and describe the changes.

- This is your detailed chain-of-thought: work through assumptions, intermediate steps, and logical derivations until the solution is reached.

Finally, provide the final answer enclosed within

\boxed{}

GROUP RELATIVE POLICY OPTIMIZATION (GRPO)

Group Relative Policy Optimization (GRPO) is a state-of-the-art Reinforcement Learning with Verifiable Rewards (RLVR) algorithm that simplifies Proximal Policy Optimization (PPO) (Schulman et al., 2017) by removing the need for a value model to estimate the baseline advantage, and has demonstrated remarkable success in enhancing the reasoning abilities of LLM. Formally, let Q denote the set of questions, $\pi_{\theta_{\text{old}}}$ be the current policy model, and $\{\mathfrak{o}_i\}_{i=1}^N$ represent a collection of Ncandidate responses sampled for a question $q \in Q$. We also define $\pi_{\theta_{\text{ref}}}$ as a fixed reference model. The training objective of GRPO is expressed as:

$$J_{\text{GRPO}}(\theta) = \mathbb{E}_{q \sim Q, \{\mathfrak{o}_i\}_{i=1}^N \sim \pi_{\theta_{\text{old}}}} \left[\frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{|\mathfrak{o}_i|} \min\left(\frac{\pi_{\theta}(\mathfrak{o}_i^t|q)}{\pi_{\theta_{\text{old}}}(\mathfrak{o}_i^t|q)} A_i, \operatorname{clip}\left(\frac{\pi_{\theta}(\mathfrak{o}_i^t|q)}{\pi_{\theta_{\text{old}}}(\mathfrak{o}_i^t|q)}, 1 - \epsilon, 1 + \epsilon \right) A_i \right) - \beta D_{\text{KL}}(\pi_{\theta} || \pi_{\text{ref}}) \right]$$
(8)

where ϵ controls the clipping range and β weights the KL regularization term. The normalized advantage A_i assigned to each response \mathfrak{o}_i is computed from group-based rewards:

$$A_i = \frac{r_i - \mu}{\sigma}, \quad \text{with } \mu = \frac{1}{N} \sum_{j=1}^{N} r_j, \quad \sigma = \sqrt{\frac{1}{N} \sum_{j=1}^{N} (r_j - \mu)^2},$$
 (9)

where $\{r_1, r_2, \dots, r_N\}$ are the scalar rewards associated with the response group $\{\mathfrak{o}_i\}_{i=1}^N$.

In GRPO, each response $\mathfrak{o} \in {\{\mathfrak{o}_i\}_{i=1}^N}$ consists of a CoT c together with its final answer. As noted in Section 2.1, token-level MDPs lack global planning and often yield redundant steps, while GRPO rewards r corresponding to o focus only on final answer correctness, overlooking reasoning quality and enabling reward hacking through superficial or verbose CoTs.

B.4 THEORETICAL PROOF

Proof. Following the framework of (Qian et al., 2025), for a fixed (q, t), the conditional error rate is

$$p_e(q, t) = 1 - \max_{y'} \Pr(y = y' \mid q, t).$$

For the binary distribution (p, 1 - p), it is known that

$$\min(p, 1 - p) \le \frac{1}{2} H_b(p),$$

where $H_b(p) = -p \log p - (1-p) \log (1-p)$ is the binary entropy. This can be generalized to the m-class case.

Lemma B.1. Let (p_1, \ldots, p_m) be a probability distribution, and let $p_{\max} = \max_i p_i$. Then

$$1 - p_{\max} \le \frac{1}{2} H(p_1, \dots, p_m).$$

Proof by induction. Base case m=2. This is exactly the binary inequality.

Induction step. Suppose the inequality holds for (m-1) classes. Consider an m-class distribution with maximum element p_1 , and let $s=1-p_1$. Merge the last two categories into one, obtaining an (m-1)-class distribution \tilde{p} . By the grouping property of Shannon entropy,

$$H(p_1,\ldots,p_{m-2},p_{m-1},p_m) = H(\tilde{p}) + (p_{m-1}+p_m)H_b\left(\frac{p_{m-1}}{p_{m-1}+p_m}\right) \ge H(\tilde{p}).$$

By the induction hypothesis,

$$s = 1 - p_1 \le \frac{1}{2}H(\tilde{p}) \le \frac{1}{2}H(p_1, \dots, p_m).$$

Thus the lemma holds for all m.

For the conditional distribution $Pr(y \mid q, t)$, the lemma implies

$$p_e(q,t) \leq \frac{1}{2}H(y \mid q,t).$$

Taking expectation over (q, t),

$$p_{\text{error}} = \mathbb{E}_{q,t}[p_e(q,t)] \le \frac{1}{2}H(y \mid q,t).$$

By the chain rule,

$$I(\hat{y}; y \mid q, t) = H(y \mid q, t) - H(y \mid \hat{y}, q, t),$$

which implies

$$H(y \mid q, t) \ge H(y \mid \hat{y}, q, t).$$

Also,

$$H(y \mid q, t) = H(y) - I(y; q, t).$$

Combining these gives

$$p_{\text{error}} \leq \frac{1}{2}H(y \mid q, t) \leq \frac{1}{2}[H(y) - I(y; \hat{y} \mid q, t)].$$

The theorem is proved.

B.5 EXPERIMENTAL PARAMETER SETUP

We conducted all experiments on eight H200 GPUs. In the supervised fine-tuning (SFT) stage, we trained Qwen2.5-7B-Instruct for 3 epochs. In the reinforcement learning (RL) stage, we adopted the GRPO algorithm, with a global batch size of 128 and a micro batch size of 4 per GPU. During rollout, the model generated 12 samples per step, including 3 analytic plans, each corresponding to 3 Chain-of-Thought (CoT) reasoning trajectories. For generation, we set temperature = 1.0 and top-p = 1.0, while for validation we used temperature = 0.6, top-p = 0.95, and n = 4. The number of RL training steps was configured as follows: LLaMA3.2-3B and Qwen2.5-7B-Instruct were trained for 350 steps, Qwen3-8B for 150 steps, and Qwen3-14B for 50 steps, with other hyperparameters kept the same across models. In addition, the learning rate (lr) was set to 1.0×10^{-6} , the weight decay (weight_decay) was 1.0×10^{-2} , the optimizer was adamw (choices: adamw or adamw_bf16), the learning-rate warmup ratio (lr_warmup_ratio) was 0. For all Qwen3-8b, max token is 4.5k and for Qwen2.5-7B-Instruct, the max token is 3.5K.