

REINFORCEMENT LEARNING WITH VERIFIABLE REWARDS IMPLICITLY INCENTIVIZES CORRECT REASONING IN BASE LLMs

Xumeng Wen^{*1}, Zihan Liu^{*†2}, Shun Zheng^{*‡1}, Zhijian Xu^{†3}, Shengyu Ye^{†1}, Zhirong Wu¹,
Xiao Liang^{†4}, Yang Wang¹, Junjie Li¹, Ziming Miao¹, Jiang Bian¹, Mao Yang¹

¹Microsoft Research Asia ²Peking University

³The Chinese University of Hong Kong ⁴University of California, Los Angeles

ABSTRACT

Reinforcement Learning with Verifiable Rewards (RLVR) has emerged as a promising paradigm for advancing the reasoning capabilities of Large Language Models (LLMs). However, a critical paradox clouds its efficacy: RLVR-tuned models often underperform their base models on the $Pass@K$ metric for solution-finding, leading to the hypothesis that RLVR merely re-weights existing reasoning paths at the cost of reasoning diversity. In this work, we resolve this contradiction by identifying the source of the problem: the $Pass@K$ metric itself is a flawed measure of reasoning, as it credits correct final answers that probably arise from inaccurate or incomplete chains of thought (CoTs). To address this, we introduce a more precise evaluation metric, $CoT-Pass@K$, which mandates that both the reasoning path and the final answer be correct. We provide a new theoretical foundation that formalizes how RLVR, unlike traditional RL, is uniquely structured to incentivize logical integrity. Our empirical results are supportive: using $CoT-Pass@K$, we observe that RLVR can incentivize the generalization of correct reasoning for all values of K . Furthermore, by analyzing the training dynamics, we find that this enhanced reasoning capability emerges early in the training process and smoothly generalizes. Our work provides a clear perspective on the role of RLVR, offers a more reliable method for its evaluation, and confirms its potential to genuinely advance machine reasoning.

1 INTRODUCTION

The successful replication of long chain-of-thought (CoT) reasoning, similar to that in OpenAI’s o1 (OpenAI, 2024), by DeepSeek-R1 (Guo et al., 2025) using the Group Relative Policy Optimization (GRPO) algorithm (Shao et al., 2024), has sparked a surge of interest within the open research community. This interest is focused on understanding, reproducing, and extending DeepSeek’s approach, as evidenced by a multitude of recent studies (Liu et al., 2025b; Hu et al., 2025; Zeng et al., 2025; Yu et al., 2025; He et al., 2025; Wen et al., 2025; Chen et al., 2025c). Fundamentally, this emerging paradigm is a form of Reinforcement Learning with Verifiable Rewards (RLVR) (Lambert et al., 2024; Guo et al., 2025; Yue et al., 2025), where a Large Language Model (LLM) acts as a policy, generating a CoT as a sequence of actions and receiving feedback on answer correctness from deterministic verifiers. This paradigm holds the promise of endowing LLMs with the ability to learn from experience through free exploration, potentially leading to unlimited intelligence (OpenAI, 2024; Guo et al., 2025; Silver & Sutton, 2025).

However, emerging concerns question the true effectiveness of RLVR. These concerns are motivated by the observation that while RLVR improves the $Pass@1$ metric, it often fails to enhance the $Pass@K$ metric compared to the base model. This phenomenon was first noted by Shao et al. (2024) during the development of GRPO. Subsequently, a systematic study by Yue et al. (2025) on

^{*}These authors contributed equally: Xumeng Wen, Zihan Liu, Shun Zheng.

[†]Work done during the internship at Microsoft Research Asia.

[‡]Correspondence to shun.zheng@microsoft.com.

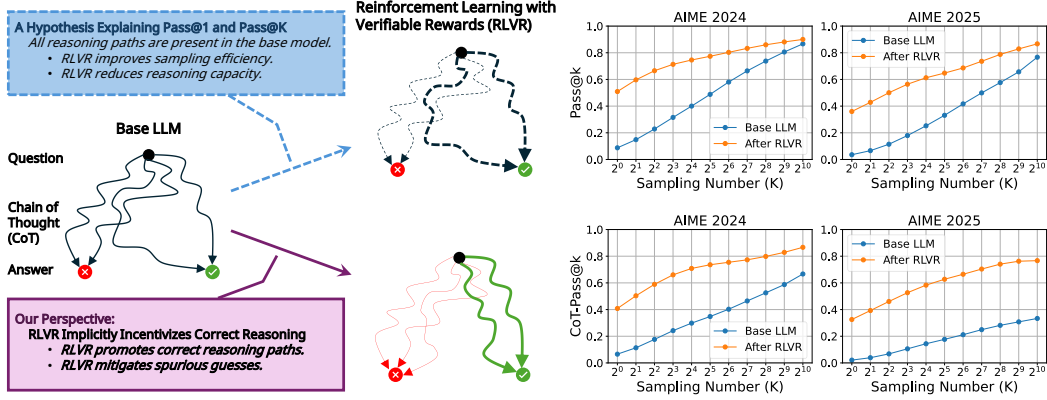


Figure 1: An illustration of our perspective: **RLVR implicitly incentivizes correct reasoning in base LLMs**. We visualize how different explanation frameworks lead to varying reasoning paths being activated, with our perspective shown in the lower left and a recent popular hypothesis explaining $Pass@K$ observations (Yue et al., 2025) summarized in the upper left. In this diagram, the line width represents the sampling probability of a reasoning path, while the color distinguishes correct paths (green) from incorrect ones (red). If all reasoning paths after applying RLVR are already present in the base model, the reasoning model merely adjusts the sampling probabilities of these existing paths (visualized in dashed lines). This hypothesis effectively accounts for the key observation shown in the upper-right part, where, for a moderately large K , a base LLM can catch up to the reasoning model after RLVR using the $Pass@K$ metric. However, we challenge this hypothesis using a refined metric, $CoT-Pass@K$, which emphasizes both the correctness of answers and the validity of the reasoning chains. As demonstrated by experimental evidence in the lower right, we observe that RLVR improves $CoT-Pass@K$ for all values of K , suggesting that **the direction of incentivization is improving the CoT correctness from $Pass@K$ of a very large K** .

various open-weight RLVR models confirmed that the $Pass@K$ metric of the base model increases at a much faster rate than its RLVR-tuned counterpart. Consequently, for a moderately large K , the base model eventually matches and surpasses the reasoning model. This led to their adventurous hypothesis: all correct reasoning paths are already present in the base model, and RLVR merely improves sampling efficiency at the cost of reducing overall reasoning capacity.

While this hypothesis has gained significant support (Zhu et al., 2025; Zhang et al., 2025; Wang et al., 2025a; Chen et al., 2025a), conflicting observations have also been reported. For instance, Liu et al. (2025a) detected the emergence of new reasoning patterns after RLVR, while they also acknowledged a loss in reasoning capacity as measured by $Pass@K$. Chen et al. (2025c) reported statistically significant improvements in $Pass@K$ for values of K up to 1024. Shojaei et al. (2025) observed similar $Pass@K$ observations on math datasets but found different patterns on puzzles with high complexity. To the best of our knowledge, no systematic explanation exists to reconcile these contradictory findings, leaving a critical question unanswered: *should we accept the hypothesis as a fundamental limitation or should we trust empirical observations that challenge the hypothesis?* In essence, we return to the core problem posed by Yue et al. (2025) and rephrase it as:

“Does RLVR genuinely incentivize new reasoning in base LLMs, and if so, why does it often fail to improve their $Pass@K$ performance?”

In this work, we propose a new perspective to resolve this debate: RLVR’s primary role is to implicitly incentivize *correct* reasoning in base LLMs, not just to find correct final answers. We argue that $Pass@K$ is an unreliable metric for evaluating true reasoning progress, as base LLMs often produce inaccurate or incomplete CoTs that coincidentally arrive at the correct solution due to their strong likelihood maximization capabilities. Under this view, the failure of RLVR to improve $Pass@K$ does not signify a failure to enhance reasoning, but rather a failure of the metric itself to capture the underlying improvement in reasoning quality.

To properly measure this phenomenon, we introduce a new metric, $CoT-Pass@K$, which evaluates success only when both the final answer and the intermediate reasoning CoT are correct. Moreover, we establish a theoretical foundation for our perspective, formalizing how RLVR’s optimization process, particularly under GRPO-style algorithms, differs from traditional RL by prioritizing the logical integrity of the reasoning path. Our theory not only aligns with our empirical results using $CoT-Pass@K$ but also explains several previously elusive phenomena observed in models like DeepSeek-R1 Guo et al. (2025).

We conduct extensive empirical validation to support our claims, but manually verifying CoT correctness at scale is challenging, especially for complex math benchmarks. We overcome this by employing a powerful yet lightweight model (DeepSeek-R1-0528-Qwen3-8B (DeepSeek, 2025)) as an automated verifier in an **LLM-as-a-CoT-Judge** paradigm, a method whose reliability we confirm through manual checks. Using this verifier, we re-evaluate the performance of a post-RLVR model (DAPO-Qwen-32B (Yu et al., 2025)) against its base model (Qwen2.5-32B-Base (Qwen, 2024)). As summarized in Figure 1, the $CoT-Pass@K$ metric clearly demonstrates that RLVR robustly incentivizes correct reasoning paths across all tested values of K (up to 1024).

Furthermore, we investigate the training dynamics to understand when this improved reasoning emerges. By reproducing GRPO-style training using the open-source DAPO recipe (Yu et al., 2025) and analyzing checkpoints, we find that RLVR begins to incentivize correct reasoning from the very early stages of training, and this capability successfully generalizes to unseen test questions. The results of our training analysis align well with our theorem, which states the implicit incentivization of correct reasoning CoTs.

The remainder of the paper is organized as follows. Section 3 presents the theoretical foundation of RLVR for LLMs. Section 4 provides empirical validation on standard benchmarks, and Section 5 analyzes the training dynamics of RLVR. Section 6 discusses limitations and future directions, Section 2 reviews related work, and Section 7 concludes the paper. Our key contributions are:

- **A New Perspective and Metric for RLVR:** We reinterpret the effect of RLVR as incentivizing correct reasoning and propose $CoT-Pass@K$ as a reliable measure. This new view addresses emerging concerns about RLVR’s efficacy and highlights its true potential.
- **A Theoretical Foundation:** We establish a theoretical foundation that distinguishes RLVR for LLMs from traditional RL for generic models by emphasizing CoT correctness. This framework formalizes the optimization dynamics of RLVR, explains previously unclear empirical results, and guides future research.
- **Empirical Validation and Training Analysis:** We observe that RLVR can improve $CoT-Pass@K$ of base LLMs for all values of K , indicating the incentivization of correct reasoning. Moreover, we observe that RLVR consistently promotes correct reasoning from early training stages and that this ability generalizes.

2 RELATED WORK

RLVR Since the release of DeepSeek-R1 (Guo et al., 2025), there has been a surge of research interest in the RLVR paradigm (Luo et al., 2025b; Liu et al., 2025b; Hu et al., 2025; Cui et al., 2025; Xie et al., 2025; Zeng et al., 2025; Yu et al., 2025; Luo et al., 2025a; Chen et al., 2025a; He et al., 2025; Wen et al., 2025; Cao et al., 2025; Liu et al., 2025a; Chen et al., 2025c). Due to the high computational cost of RLVR, most studies have focused on small- to medium-sized models (up to 32B parameters). These studies span a wide range of aspects, including training data curation, objective design, hyperparameter tuning, base model selection, and various insightful observations. However, only a few studies have addressed the theoretical foundations of RLVR. In this work, we argue that RLVR for LLMs should be understood from a different perspective—one that emphasizes the correctness of reasoning paths. We hope our theoretical perspective and empirical findings will inspire the community to develop more efficient and effective RLVR approaches, unlocking its broader potential across diverse applications.

Debates on Whether RLVR Really Incentivizes Since Yue et al. (2025) raised the insightful question of whether RLVR truly incentivizes improvements beyond the base LLMs, and conducted extensive empirical experiments to demonstrate the wide applicability of their key hypothesis—that

RLVR does not improve $Pass@K$ for the base LLM because all reasoning paths are already present in the base model—there have been varying perspectives on this hypothesis. Some researchers agree with this viewpoint (Wang et al., 2025b; Zhu et al., 2025; Zhang et al., 2025; Wang et al., 2025a; Chen et al., 2025a), while others report contradictory findings (Liu et al., 2025a; Chen et al., 2025c; Shojaei et al., 2025), as discussed in the introduction. There is currently no fundamental understanding to resolve these debates. Liu et al. (2025a) speculated that previous RLVR experiments may have been conducted within a single domain (e.g., math) and were optimized for limited gradient steps before true exploration could occur. Shojaei et al. (2025) suggested that the complexity of puzzles might be the key factor. Chen et al. (2025c) presented statistically significant empirical results to justify that their model indeed improves $Pass@K$, particularly highlighting a persistent gap on the LiveCodeBench v6 (Jain et al., 2025), leading them to conclude that the base model is likely guessing. In this work, we align with the intuition of Chen et al. (2025c) and believe in the rationality of their empirical results. Our findings also suggest that on challenging, live benchmarks, base LLMs struggle to guess, and their limitations in reasoning become clearly evident.

The Importance of Correct CoTs Recent studies have also highlighted the importance of verifying the correctness of CoTs (Arcuschin et al., 2025; McGinness & Baumgartner, 2025; Shojaei et al., 2025). However, their approaches focus on defining synthetic reasoning tasks where the correctness of reasoning CoTs can be verified easily. While this is an interesting and effective approach for fully examining reasoning correctness, it is difficult to apply to unstructured reasoning scenarios, such as in math and code. In this work, we argue that the LLM-as-a-CoT-Judge paradigm could play a crucial role in more general reasoning tasks, and emphasize the pressing need for the design of evaluation benchmarks to assess the reliability of emerging LLM verifiers. In the meanwhile, we note there is a contemporary study also advocating this paradigm (Jiang et al., 2025), and they mainly consider education and healthcare domains.

3 A THEORETICAL FOUNDATION OF RLVR FOR LLMs

In this section, we establish a theoretical foundation for how RLVR, as implemented in the GRPO algorithm (Shao et al., 2024), incentivizes the generation of correct reasoning CoTs, which we define as being both logically accurate and complete.

A key distinction must be made between RLVR and traditional RL. Base LLMs, owing to their powerful likelihood estimation capabilities obtained during pre-training, can generate numerous incorrect or incomplete CoTs that coincidentally arrive at a correct final answer. In contrast, traditional RL simply optimizes for action trajectories that yield high rewards, without necessarily verifying the intrinsic correctness of each action along the path. For instance, in the Go game (Silver et al., 2017), every action is valid once the simulation environment is setup correctly.

In the context of LLMs, we argue that the core principle of RLVR is fundamentally different. It is not merely about reaching a correct answer, but about exploring the immense reasoning space with broad prior knowledge and about identifying and reinforcing logically rigorous CoTs. To formalize this principle, we now elaborate on our problem formulation, key assumptions, the resulting theorem, and some discussions of its implications.

3.1 PROBLEM SETUP

For each prompt q , we sample G responses $\mathbf{Y} = \{y_1, y_2, \dots, y_G\}$ from policy π_θ . Let c_i be the CoT in response y_i , and a_i the final answer. Define correctness indicators:

$$\mathcal{I}_{\text{CoT}}(c_i) = \begin{cases} 1 & \text{if } c_i \text{ is correct (logically accurate and complete)} \\ 0 & \text{otherwise} \end{cases}, \quad (1)$$

$$\mathcal{I}_{\text{Ans}}(a_i) = \begin{cases} 1 & \text{if } a_i \text{ is correct} \\ 0 & \text{otherwise} \end{cases}. \quad (2)$$

We have a verifiable reward $R(y_i)$ that is binary and determined solely by answer correctness:

$$R(y_i) = \mathcal{I}_{\text{Ans}}(a_i). \quad (3)$$

The GRPO advantage $\hat{A}(y_i)$ is computed as:

$$\hat{A}(y_i) = \frac{R(y_i) - \mu_{\mathbf{Y}}}{\sigma_{\mathbf{Y}}}, \quad \mu_{\mathbf{Y}} = \frac{1}{G} \sum_{j=1}^G R(y_j), \quad \sigma_{\mathbf{Y}} = \sqrt{\frac{1}{G} \sum_{j=1}^G (R(y_j) - \mu_{\mathbf{Y}})^2}. \quad (4)$$

We consider a simplified GRPO gradient update:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{G} \sum_{i=1}^G \hat{A}(y_i) \nabla_{\theta} \log \pi_{\theta}(y_i | q). \quad (5)$$

3.2 THE THEOREM

Given the following assumptions, we establish Theorem 1.

- **Logical Coherence:** Compared with incorrect CoTs, correct CoTs have higher probabilities to induce correct answers since base LLMs has been pretrained over massive corpora to establish strong logical priors:

$$P(\mathcal{I}_{\text{Ans}}(a_i) = 1 | \mathcal{I}_{\text{CoT}}(c_i) = 1) = \alpha > P(\mathcal{I}_{\text{Ans}}(a_i) = 1 | \mathcal{I}_{\text{CoT}}(c_i) = 0) = \beta \quad (6)$$

- **Stable Advantage Estimation:** A sufficiently large group G to ensure statistically stable advantage estimates, and this group is learnable ($\sigma_{\mathbf{Y}} > 0$).

Theorem 1 (GRPO Implicitly Incentivizes Correct Reasoning) *Given the above problem setup and two assumptions, the expected GRPO advantage $\mathbb{E}[\hat{A}(y_i)]$ satisfies:*

$$\mathbb{E}[\hat{A}(y_i) | \mathcal{I}_{\text{CoT}}(c_i) = 1] > 0, \quad (7)$$

$$\mathbb{E}[\hat{A}(y_i) | \mathcal{I}_{\text{CoT}}(c_i) = 0] < 0, \quad (8)$$

for any prompt q , where $\hat{A}(y_i)$ is defined in equation 4. Consequently, GRPO policy updates (equation 5) increase the likelihood of generating correct CoTs.

Proof 1 Let $p_c = P(\mathcal{I}_{\text{CoT}}(c_i) = 1)$ be the current probability of generating a correct CoT. The expected reward for a response y_i is:

$$\mathbb{E}[R(y_i)] = \begin{cases} \alpha & \text{if } \mathcal{I}_{\text{CoT}}(c_i) = 1 \\ \beta & \text{if } \mathcal{I}_{\text{CoT}}(c_i) = 0 \end{cases} \quad (9)$$

The group-level expected reward $\mu \triangleq \mathbb{E}[\mu_{\mathbf{Y}}]$ is:

$$\mu = p_c \alpha + (1 - p_c) \beta. \quad (10)$$

For large G , the group mean $\mu_{\mathbf{Y}}$ and variance $\sigma_{\mathbf{Y}}^2$ concentrate around their expectations:

$$\mu_{\mathbf{Y}} \xrightarrow{G \rightarrow \infty} \mu \quad (11)$$

$$\sigma_{\mathbf{Y}}^2 \xrightarrow{G \rightarrow \infty} \sigma^2 > 0. \quad (12)$$

The expected advantage conditional on CoT correctness is:

$$\mathbb{E}[\hat{A}(y_i) | \mathcal{I}_{\text{CoT}}(c_i) = 1] \xrightarrow{G \rightarrow \infty} \frac{\alpha - \mu}{\sigma} \quad (13)$$

$$\mathbb{E}[\hat{A}(y_i) | \mathcal{I}_{\text{CoT}}(c_i) = 0] \xrightarrow{G \rightarrow \infty} \frac{\beta - \mu}{\sigma}. \quad (14)$$

Substituting equation 10 into equation 13 and equation 14:

$$\mathbb{E}[\hat{A}(y_i) | \text{correct CoT}] \rightarrow \frac{(1 - p_c)(\alpha - \beta)}{\sigma} \quad (15)$$

$$\mathbb{E}[\hat{A}(y_i) | \text{incorrect CoT}] \rightarrow \frac{-p_c(\alpha - \beta)}{\sigma}. \quad (16)$$

Since $\alpha > \beta$ (by equation 6 under the assumption of logical coherence) and $\sigma > 0$, we have:

$$\begin{aligned}(1 - p_c)(\alpha - \beta)/\sigma &> 0, \\ -p_c(\alpha - \beta)/\sigma &< 0,\end{aligned}$$

proving inequalities equation 7 and equation 8.

The GRPO policy gradient update in equation 5, $\nabla_\theta J(\theta) \approx \frac{1}{G} \sum_{i=1}^G \hat{A}(y_i) \nabla_\theta \log \pi_\theta(y_i | q)$, on average increases the likelihood of responses with $\hat{A}(y_i) > 0$ (correct CoTs) and decreases it for $\hat{A}(y_i) < 0$ (incorrect CoTs). Thus, p_c increases monotonically.

Discussions on (p_c, α, β) Theorem 1 demonstrates that GRPO inherently aligns policy updates with correct reasoning, even for base models with low initial p_c . The driving factor is the gap $\alpha - \beta > 0$, which amplifies the advantage difference between correct and incorrect CoTs. As training progresses and α increases (due to more sound reasoning) while β decreases (reducing spurious correlations), causing the gap to widen and further accelerating coherent reasoning. As $p_c \rightarrow 1$, $(\alpha - \beta)$ may approach 1 in a faster pace because generating short answers is typically much easier than producing long correct CoTs, then $\mathbb{E}[\hat{A}(y_i) | \text{correct CoT}] \rightarrow 0$, ensuring convergence.

Discussions on (μ, σ^2) From equation 10, we know that the group reward mean is given by $\mu = p_c \alpha + (1 - p_c) \beta$. Furthermore, we can derive the exact formula for the variance σ^2 in equation 12 and analyze their impacts together with p_c , α , and β on policy iterations.

The sample variance σ_Y^2 converges to the true variance σ^2 :

$$\sigma_Y^2 = \frac{1}{G} \sum_{j=1}^G (R(y_j) - \mu_Y)^2 \xrightarrow{G \rightarrow \infty} \text{Var}(R(y_j)) \equiv \sigma^2,$$

where $\text{Var}(R(y_j))$ can be computed using the law of total variance:

$$\text{Var}(R(y_j)) = \underbrace{\text{Var}(\mathbb{E}[R(y_j) | \mathcal{I}_{\text{CoT}}(c_j)])}_{\text{Variance of conditional expectation}} + \underbrace{\mathbb{E}[\text{Var}(R(y_j) | \mathcal{I}_{\text{CoT}}(c_j))]}_{\text{Expectation of conditional variance}}.$$

First term:

$$\mathbb{E}[R(y_j) | \mathcal{I}_{\text{CoT}}(c_j)] = \begin{cases} \alpha & \text{if } \mathcal{I}_{\text{CoT}}(c_j) = 1 \\ \beta & \text{if } \mathcal{I}_{\text{CoT}}(c_j) = 0 \end{cases}.$$

The random variable $\mathbb{E}[R(y_j) | \mathcal{I}_{\text{CoT}}(c_j)]$ has variance:

$$\text{Var}(\mathbb{E}[R(y_j) | \mathcal{I}_{\text{CoT}}(c_j)]) = (\alpha - \beta)^2 p_c (1 - p_c).$$

Second term:

$$\text{Var}(R(y_j) | \mathcal{I}_{\text{CoT}}(c_j)) = \begin{cases} \alpha(1 - \alpha) & \text{if } \mathcal{I}_{\text{CoT}}(c_j) = 1 \\ \beta(1 - \beta) & \text{if } \mathcal{I}_{\text{CoT}}(c_j) = 0 \end{cases},$$

so its expectation is:

$$\mathbb{E}[\text{Var}(R(y_j) | \mathcal{I}_{\text{CoT}}(c_j))] = p_c \alpha (1 - \alpha) + (1 - p_c) \beta (1 - \beta).$$

Thus:

$$\sigma^2 = (\alpha - \beta)^2 p_c (1 - p_c) + p_c \alpha (1 - \alpha) + (1 - p_c) \beta (1 - \beta). \quad (17)$$

Substituting μ and σ into equation 15 and equation 16, we have

$$\begin{aligned}\mathbb{E}[\hat{A}(y_i) | \text{correct CoT}] &\rightarrow \frac{(1 - p_c)(\alpha - \beta)}{\sqrt{(\alpha - \beta)^2 p_c (1 - p_c) + p_c \alpha (1 - \alpha) + (1 - p_c) \beta (1 - \beta)}}, \\ \mathbb{E}[\hat{A}(y_i) | \text{incorrect CoT}] &\rightarrow \frac{-p_c(\alpha - \beta)}{\sqrt{(\alpha - \beta)^2 p_c (1 - p_c) + p_c \alpha (1 - \alpha) + (1 - p_c) \beta (1 - \beta)}}.\end{aligned}$$

An ideal pre-training on a high-capacity model could help to ensure that $\alpha \rightarrow 1$ and $\beta \rightarrow 0$ at the beginning of RLVR. In this condition, we have the following advantage estimates:

$$\mathbb{E}[\hat{A}(y_i) \mid \text{correct CoT}] \rightarrow \sqrt{\frac{1-p_c}{p_c}}, \quad \mathbb{E}[\hat{A}(y_i) \mid \text{incorrect CoT}] \rightarrow -\sqrt{\frac{p_c}{1-p_c}}.$$

In this ideal scenario, the role of human would be to prepare a comprehensive and diverse set of questions and answers, leveraging RLVR to automatically incentivize the model’s reasoning capabilities. However, in practice—the “unideal case”—it is often necessary to first fine-tune the base LLM to align its output with a proper reasoning distribution before applying RLVR.

Discussions on Key Observations in RLVR Grounded in our theoretical analysis, we can now provide our unique explanations for several previously elusive yet important observations reported in DeepSeek-R1 (Guo et al., 2025).

Our Explanation of the Observation “*DeepSeek-R1-Zero achieved remarkable Pass@K performance on AIME 2024 but encountered challenges such as poor readability and language mixing.*”: Even DeepSeek-V3 (Liu et al., 2024) cannot guarantee ideal conditions where $\alpha \rightarrow 1, \beta \rightarrow 0$. As a result, cold-start data is required to rectify prior logic biases, motivating the R1 approach.

Our Explanation of the Observation “*The R1-Zero approach did not work well for the 32B dense model, yet distillation can be very effective.*”: Key factors such as (p_c, α, β) for the 32B base model are in an even worse state, causing pure RLVR to converge to suboptimal local solutions. Based on our analysis, the key to effective reasoning lies in learning correct CoTs. Therefore, the distillation approach can efficiently teach an LLM how to reason properly.

Our Explanation of the Observation “*The average response length of DeepSeek-R1-Zero naturally increases during training.*”: On average, long CoTs have higher probabilities than short CoTs to generate correct answers because more tokens can enable problem solving in finer-grained steps and may also introduce more spurious correlations. Replacing “correct v.s. incorrect” with “long v.s. short” in equation 6 leads to the conclusion that long CoTs being naturally incentivized. For simple problems, long CoTs may be regarded as an improper model bias, which could be the root cause of widely observed “over-thinking” phenomena (Chen et al., 2025b).

Discussions on Exceptional Cases We acknowledge that the assumption of logical coherence (equation 6) may not always hold, potentially leading to the reinforcement of incorrect CoTs. As previously discussed, base LLMs may retain inherent biases from pre-training—though incorrect, these biases might coincidentally yield the right final answer due to spurious correlations. In such cases, improper model biases could be unintentionally reinforced. Consequently, we believe that additional techniques, such as learning from human feedback (Ouyang et al., 2022) or off-policy-guided learning (Yan et al., 2025), may prove essential in addressing these misalignments.

3.3 KEY METRICS TO MEASURE

For each prompt q with G responses, we define the number of correct answers and the number of correct CoTs (with correct final answers) as:

$$C = \sum_{i=1}^G \mathcal{I}_{\text{Ans}}(a_i) \quad (\text{Number of correct answers}) \quad (18)$$

$$D = \sum_{i=1}^G \mathcal{I}_{\text{CoT}}(c_i) \cdot \mathcal{I}_{\text{Ans}}(a_i) \quad (\text{Correct CoTs with correct answers}) \quad (19)$$

We estimate $\text{Pass}@K$ using the method introduced by Chen et al. (2021); Yue et al. (2025). Accordingly, we define the specific calculations for per-prompt key metrics for any $K \leq G$ as:

$$Pass@K^{(q)} = 1 - \frac{\binom{G-C}{K}}{\binom{G}{K}} \quad (\text{Prob. of at least one correct answer}) \quad (20)$$

$$CoT-Pass@K^{(q)} = 1 - \frac{\binom{G-D}{K}}{\binom{G}{K}} \quad (\text{Prob. of at least one correct CoT \& answer}) \quad (21)$$

$$P(CA)^{(q)} = \frac{C}{G} \quad (\text{Fraction of correct answers} = Pass@1^{(q)}) \quad (22)$$

$$P(CC|CA)^{(q)} = \frac{D}{C} \quad (\text{Fraction of correct CoTs in correct answers}) \quad (23)$$

The overall (averaged) metrics across M prompts are given by:

$$Pass@K = \frac{1}{M} \sum_{q=1}^M Pass@K^{(q)} \quad (24)$$

$$CoT-Pass@K = \frac{1}{M} \sum_{q=1}^M CoT-Pass@K^{(q)} \quad (25)$$

$$P(CA) = \frac{1}{M} \sum_{q=1}^M P(CA)^{(q)} \quad (26)$$

$$P(CC|CA) = \frac{1}{M} \sum_{q=1}^M P(CC|CA)^{(q)} \quad (27)$$

4 REVISITING PASS@K EXPERIMENTS WITH CoT-PASS@K

We revisit the $Pass@K$ experiments on popular math benchmarks using EvalHub (Ye, 2025), introducing $CoT-Pass@K$ to provide a more accurate assessment of reasoning. A prominent challenge in this analysis is the verification of massive volumes of long and complex CoTs, a task that requires expert-level mathematical knowledge and is prohibitively difficult to perform manually at scale. To address this, we leverage the recently released DeepSeek-R1-0528 (DeepSeek, 2025), employing its distilled 8B variant, DeepSeek-R1-0528-Qwen3-8B, as a powerful yet lightweight verifier. We developed a specific prompt template for this task (see Appendix A.4). Following automatic verifications at scale, we confirmed the reliability of this LLM-as-a-CoT-Judge paradigm by manually verifying its judgments on some of the most difficult problems (see Appendix A.5).

To mitigate potential errors from the LLM verifier, which is powerful but not infallible, we verify each CoT multiple times. We then determine the final CoT correctness using three distinct strategies to ensure the robustness of our findings: *any-correct* (at least one verification returns correct), *all-correct* (all verifications must return correct), and *majority-correct* (a majority vote determines the outcome). In Appendix A.3, we have justified that this multi-verification system can mitigate both false positives and false negatives.

Figure 2 presents a comparison between the base LLM and its post-RLVR counterpart using both $Pass@K$ and $CoT-Pass@K$. The $Pass@K$ results (top row) confirm the observations in (Yue et al., 2025): the performance of the base LLM appears to catch up and even surpass the post-RLVR model as K increases. However, in stark contrast, the $CoT-Pass@K$ results on AIME 2024 and AIME 2025 reveal a persistent and significant performance gap between the models across all values of K (up to 1024). This gap is especially pronounced on AIME 2025, as it is free from data contamination, having been released after the base model’s training cutoff. Manual inspection of numerous cases confirms that the base LLM frequently arrives at correct answers through flawed reasoning (see examples in Appendix A.5.1 and A.5.2). These flawed solutions, which inflate the standard $Pass@K$ score, are correctly filtered out by our $CoT-Pass@K$ metric. Conversely, the post-RLVR model consistently produces rigorous reasoning chains, as evidenced by its high scores even under the strict *all-correct* verification strategy.

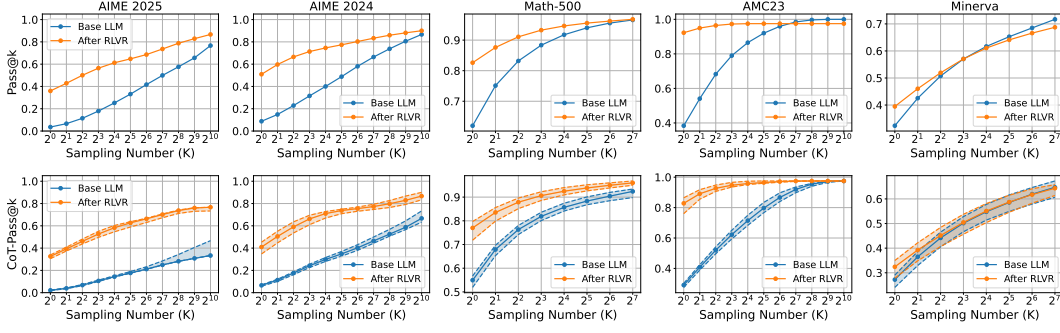


Figure 2: Comparisons of $Pass@K$ (the top row) and $CoT-Pass@K$ (the bottom row) on five math benchmarks (different columns) to show different observations of how RLVR incentivizes reasoning in LLMs. Here the base LLM is Qwen2.5-32B, and the post-RLVR model is DAPO-Qwen-32B. For $CoT-Pass@K$, we perform multiple verifications for each CoT using DeepSeek-R1-0528-Qwen3-8B, and display the results determined by *any-correct*, *all-correct*, and *majority-correct* strategies, which constitute the shaded area in lower subplots.

Nevertheless, we observe that on other benchmarks like Math-500 and AMC23, the incentivizing effects of RLVR are less apparent, as the base LLM is already capable of solving these problems correctly within a few trials. This could be because 1) the problems are simple enough for the base LLM to solve using its existing knowledge, or 2) the problems were part of its pre-training data so the base LLM can easily recall a correct solution given multiple trials. It is difficult to distinguish these possibilities without knowing the training data recipe of Qwen2.5-32B. Furthermore, on the Minerva benchmark, the post-RLVR model shows no improvement. This is likely attributable to a domain mismatch, as Minerva contains many physics problems and more free-form answers, whereas the DAPO training data was restricted to math problems formatted to produce integer answers.

Our theoretical framework ensures that RLVR incentivizes correct reasoning for training prompts, but it does not guarantee generalization across all scenarios. Therefore, the observed evaluation variations do not challenge the validity of our framework. The results on AIME 2024 and AIME 2025 already demonstrate the generalization of correctly reasoned generations incentivized during training. Moreover, these differing generalization behaviors highlight the critical importance of evaluating RLVR on challenging, contamination-free benchmarks to accurately assess its impact on model reasoning capabilities. They also underscore the need for curating comprehensive and diverse datasets to effectively scale RLVR, as demonstrated in (Liu et al., 2025a; Chen et al., 2025c).

5 ANALYZING THE TRAINING DYNAMICS OF RLVR

The existence of generalizable, incentivized correct reasoning on AIME 2024 and AIME 2025 motivates us to investigate when such incentivization emerges during RLVR training. To this end, we adopt the open-sourced DAPO training recipe (Yu et al., 2025), which follows the R1-zero approach starting from the base LLM Qwen2.5-32B and claims to achieve results better than DeepSeek-R1 (Guo et al., 2025) on the same base model.

Our reproduction was conducted on 32 AMD MI300X GPUs using the VERL framework (Sheng et al., 2025), and ran for over two weeks. While our run did not reproduce the $Pass@1$ accuracy above 50% as reported by Yu et al. (2025), we reached a comparable performance of around 44% $Pass@1$, in line with a third-party reproduction (Chen et al., 2025a). We use the same verifier introduced in Section 4 to assess the correctness of both training and evaluation rollouts.

Figure 3 summarizes the training dynamics of our DAPO reproduction. We observe that RLVR begins to incentivize correct reasoning from the very beginning, as evidenced by increased $P(CC|CA)^{(q)}$ values in the early training steps shown in Figures 3(a) and 3(b). These incentivized reasoning capabilities translate into improved generalization on unseen questions, as demonstrated by notable gains in $CoT-Pass@K$ on AIME 2024 within the first 20 training steps in Figure 3(c). Note that each training step here corresponds to one round of PPO-style optimization (Schulman

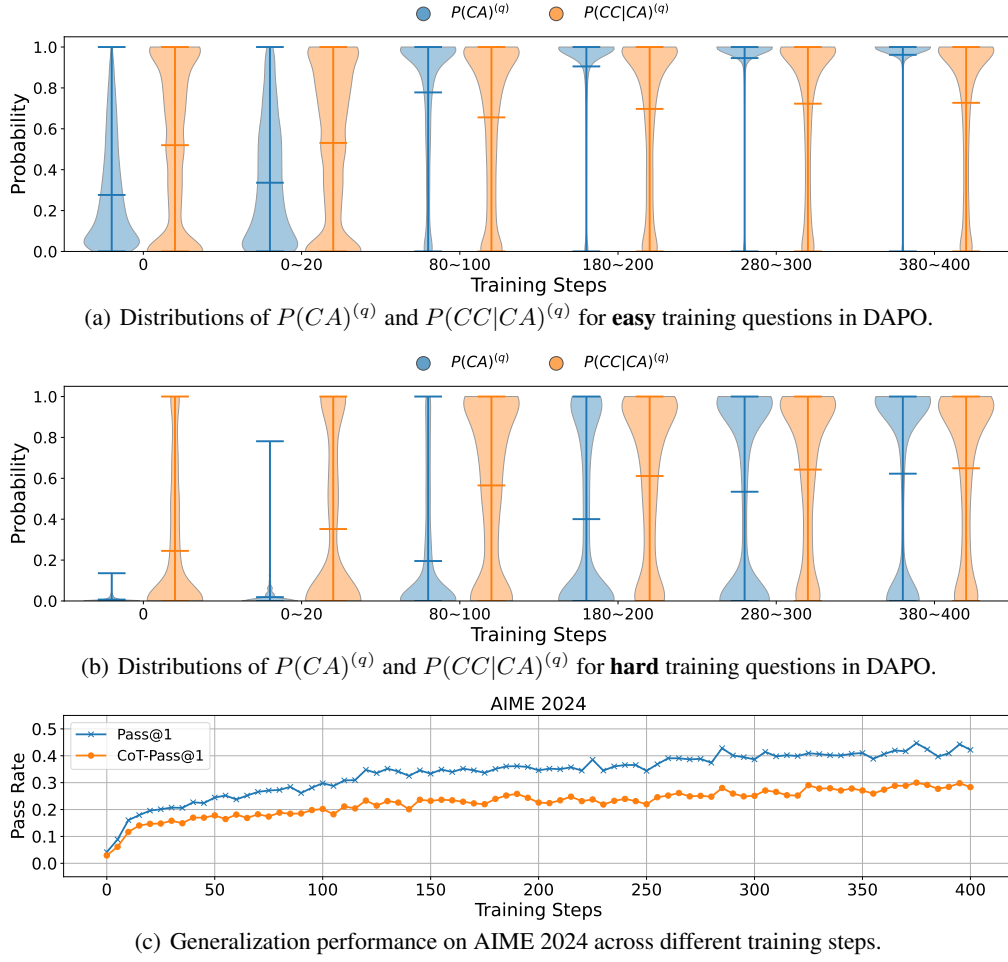


Figure 3: Training dynamics of our reproduction of DAPO, following the recipes provided by Yu et al. (2025). The top two subfigures use violin plots to visualize the evolution of $P(CA)^{(q)}$ (the fraction of correct answers for prompt q) and $P(CC|CA)^{(q)}$ (the fraction of correct CoTs within the correct answers for prompt q) over the course of DAPO training. Subfigure (a) shows results for the easy training questions, while (b) presents the corresponding distributions for the hard questions. We divide training questions into easy and hard by sampling 64 rollouts for each of the 17k training questions using Qwen2.5-32B, and labeling those with at least one correct answer as easy. Subfigure (c) presents the generalization performance on AIME 2024 across different training steps.

et al., 2017), which includes 16 gradient updates, according to the DAPO training script. Thus, we see that correct reasoning abilities begin to generalize after only a few gradient updates.

Furthermore, the incentivization of correct reasoning on training questions appears to be a continuous process, as reflected by the steady increase in the mean of $P(CC|CA)^{(q)}$ throughout training, for both easy and hard questions. Meanwhile, we again observe that $P(CA)^{(q)}$ (equivalent to $Pass@1^{(q)}$) is an unreliable metric, particularly for easy training questions. As shown in Figure 3(a), the distribution of $P(CA)^{(q)}$ becomes highly skewed toward 1.0 after 180 steps, misleadingly suggesting that most questions are perfectly solved. However, examining the distribution of $P(CC|CA)^{(q)}$ reveals that a substantial fraction of responses still contain flawed reasoning. We suspect this is one of the reasons behind the difficulty of achieving strong results with Qwen2.5-32B using the R1-zero approach.

In addition, for hard questions, we observe that the mean of $P(CA)^{(q)}$ increases more quickly than that of $P(CC|CA)^{(q)}$, albeit at a slower rate compared to the easy-question setting. In both

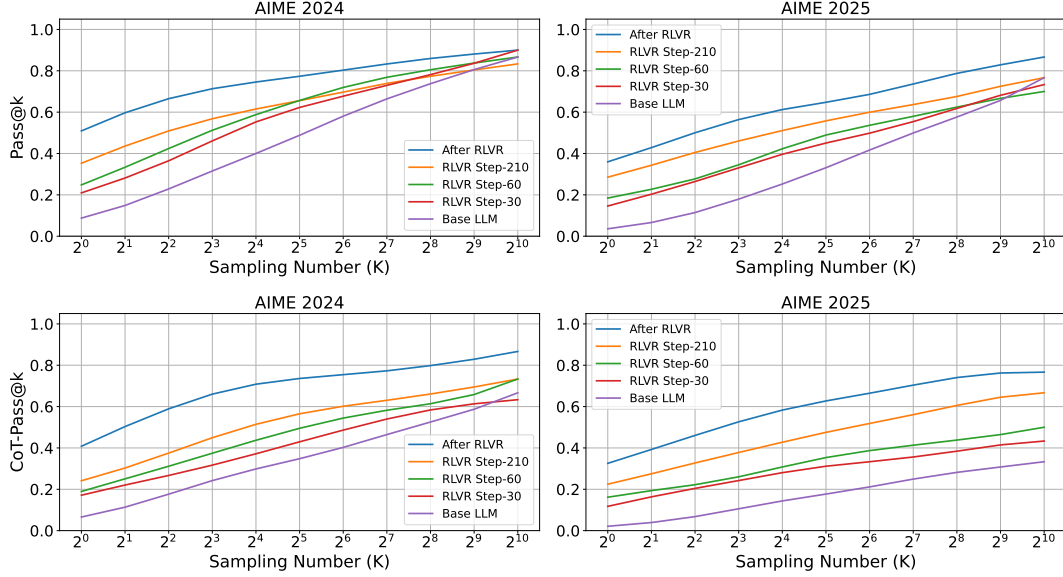


Figure 4: Revisiting $Pass@K$ and $CoT-Pass@K$ experiments on AIME 2024 and AIME 2025 using early and mid-stage checkpoints of our DAPO reproduction. The base LLM and post-RLVR model are Qwen2.5-32B and DAPO-Qwen-32B, respectively.

cases, improving $P(CC|CA)^{(q)}$ proves to be a slow and challenging process. Since our analysis shows that incentivizing correct CoTs is key to improving reasoning capabilities, we believe that future research should explore novel mechanisms to accelerate the improvement of $P(CC|CA)^{(q)}$, thereby enhancing both the efficiency and effectiveness of RLVR.

To further support the claim that RLVR incentivizes correct reasoning from the start in a smooth and consistent manner, we conduct additional evaluations using early and mid-stage checkpoints from our DAPO reproduction. Figure 4 presents the corresponding $Pass@K$ and $CoT-Pass@K$ results on AIME 2024 and AIME 2025 with K scaled up to 1024, while the initial DAPO experiment adopts $K = 16$ in training and $K = 32$ in testing. These results more clearly reveal the underlying incentivization of reasoning as captured by $CoT-Pass@K$. The contamination-free AIME 2025 benchmark provides especially clear evidence of this effect across all tested values of K .

We believe these empirical findings from the training dynamics of RLVR strongly validate the theoretical framework proposed in this work.

6 DISCUSSIONS

Limitations A key limitation of our study lies in the use of a LLM as the verifier for the correctness of reasoning CoTs, due to the prohibitive cost of manually checking a large volume of generated reasoning paths. To mitigate this, we present extensive case studies in Appendix A.5 to demonstrate that DeepSeek-R1-0528-Qwen3-8B functions as a relatively robust verifier across multiple math benchmarks. Furthermore, we apply multiple verification calls to obtain $CoT-Pass@K$ metrics under various criteria, including *any-correct*, *majority-correct*, and *all-correct*, in order to balance between false positives and false negatives. Another limitation is the current focus on math reasoning and a limited number of post-RLVR models. We plan to broaden the scope in future work by incorporating more reasoning domains and more models.

Call for Live, Challenging Benchmarks Static benchmarks developed prior to the release of modern base models are increasingly susceptible to contamination risks, potentially undermining the reliability of observed improvements. In response, we emphasize the need for *live benchmarks* that evolve over time, as suggested in recent studies (Jain et al., 2025; White et al., 2025). Additionally,

we agree with the viewpoint of Yao (2025) that future research advancements may rely more on designing new evaluations, benchmarks, and environments.

Call for Lightweight yet Powerful CoT Verifiers While DeepSeek-R1-0528-Qwen3-8B serves as a useful CoT verifier, it is not infallible. Conflicting verification results across multiple queries reveal the challenges of false-positive and false-negative verifications. To tackle this, we combine multiple verification strategies, including different voting rules, to improve robustness. Looking forward, there is a pressing need for light yet reliable CoT verifiers that can serve as standardized evaluators beyond the coarse-grained Pass@K metric. This direction also relates to previous studies on process reward modeling (Lightman et al., 2024; Uesato et al., 2022; Wang et al., 2024).

Scaling RLVR or Scaling Pre-Training While the scaling of pre-training has led to transformative progress in LLMs (Kaplan et al., 2020; Liu et al., 2024), enabling the transition to the era of artificial general intelligence, we argue that scaling RLVR could be equally pivotal, given the empirical evidences and theoretical foundation that all demonstrate its real incentivization beyond base LLMs. As modern LLMs approach the limits of language token exposure, learning from experience (Silver & Sutton, 2025) may represent the next leap. Recent efforts by leading research teams suggest a growing emphasis on this direction (Guo et al., 2025; DeepSeek, 2025; Gemini, 2024; Grok, 2025; OpenAI, 2025; Qwen, 2025; Gemini, 2025; Anthropic, 2025; Mistral.AI, 2025). For the broad open research community, understanding the foundations and limitations of current RLVR algorithms is crucial to push this direction further.

New RLVR Algorithms and Beyond With our insight that RLVR implicitly incentivizes correct reasoning in base LLMs, we anticipate the development of new algorithmic paradigms. These may include optimization formulations or objective functions, such as policy-gradient approaches (Sutton et al., 1999; Schulman et al., 2017), new likelihood-based optimization objectives (Chen et al., 2025a; Zhu et al., 2025), and preference optimization frameworks (Rafailov et al., 2023; Su et al., 2025). The key principle is that the new algorithms should be designed to more directly incentivize correct reasoning paths, alleviating inherent logical biases in base LLMs.

7 CONCLUSION

In this work, we revisit the fundamental question of whether RLVR genuinely incentivizes new reasoning in base LLMs. Through a new perspective, we propose that RLVR implicitly promotes correct CoT reasoning, and we introduce the CoT-Pass@K metric to better capture this effect. Our theoretical analysis explains the underlying optimization mechanisms distinguishing RLVR from traditional RL. And we conduct extensive empirical validation and training analysis to further support the new perspective and our theory.

These findings not only resolve conflicting observations in prior work but also illuminate the untapped potential of RLVR in aligning LLMs with human reasoning systems. We envision a promising future where RLVR serves as a cornerstone for developing LLMs that learn through interaction, self-correction, and verifiable reasoning.

REFERENCES

- Team Anthropic. Introducing Claude 4. <https://www.anthropic.com/news/claude-4>, 2025. [Released 23-05-2025].
- Iván Arcuschin, Jett Janiak, Robert Krzyzanowski, Senthooan Rajamanoharan, Neel Nanda, and Arthur Conmy. Chain-of-thought reasoning in the wild is not always faithful. *arXiv preprint arXiv:2503.08679*, 2025.
- Shiyi Cao, Sumanth Hegde, Dacheng Li, Tyler Griggs, Shu Liu, Eric Tang, Jiayi Pan, Xingyao Wang, Akshay Malik, Graham Neubig, Kourosh Hakhmaneshi, Richard Liaw, Philipp Moritz, Matei Zaharia, Joseph E. Gonzalez, and Ion Stoica. SkyRL-v0: Train real-world long-horizon agents via reinforcement learning. <https://novasky-ai.notion.site/skyrl-v0>, 2025.

- Huayu Chen, Kaiwen Zheng, Qinsheng Zhang, Ganqu Cui, Yin Cui, Haotian Ye, Tsung-Yi Lin, Ming-Yu Liu, Jun Zhu, and Haoxiang Wang. Bridging supervised learning and reinforcement learning in math reasoning. *arXiv preprint arXiv:2505.18116*, 2025a.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wanxiang Che. Towards reasoning era: A survey of long chain-of-thought for reasoning large language models. *arXiv preprint arXiv:2503.09567*, 2025b.
- Yang Chen, Zhuolin Yang, Zihan Liu, Chankyu Lee, Peng Xu, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. AceReason-Nemotron: Advancing math and code reasoning through reinforcement learning. *arXiv preprint arXiv:2505.16400*, 2025c.
- 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.
- Team DeepSeek. DeepSeek-R1-0528 release. <https://api-docs.deepseek.com/news/news250528>, 2025. [Released 28-05-2025].
- Team Gemini. Introducing Gemini 2.0: our new AI model for the agentic era. <https://blog.google/technology/google-deepmind/google-gemini-ai-update-december-2024/>, 2024. [Released 11-12-2024].
- Team Gemini. Gemini 2.5: Our most intelligent models are getting even better. <https://blog.google/technology/google-deepmind/google-gemini-updates-io-2025/>, 2025. [Released 20-05-2025].
- Team Grok. Grok 3 Beta — The age of reasoning agents. <https://x.ai/news/grok-3>, 2025. [Released 19-02-2025].
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Jujie He, Jiakai Liu, Chris Yuhao Liu, Rui Yan, Chaojie Wang, Peng Cheng, Xiaoyu Zhang, Fuxiang Zhang, Jiacheng Xu, Wei Shen, et al. Skywork open reasoner 1 technical report. *arXiv preprint arXiv:2505.22312*, 2025.
- Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xiangyu Zhang, and Heung-Yeung Shum. Open-Reasoner-Zero: An open source approach to scaling up reinforcement learning on the base model. *arXiv preprint arXiv:2503.24290*, 2025.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. LiveCodeBench: Holistic and contamination free evaluation of large language models for code. In *ICLR*, 2025.
- Enyi Jiang, Changming Xu, Nischay Singh, and Gagandeep Singh. Misaligning reasoning with answers—a framework for assessing LLM CoT robustness. *arXiv preprint arXiv:2505.17406*, 2025.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *ICLR*, 2024.

- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. DeepSeek-V3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong. ProRL: Prolonged reinforcement learning expands reasoning boundaries in large language models. *arXiv preprint arXiv:2505.24864*, 2025a.
- 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*, 2025b.
- Michael Luo, Sijun Tan, Roy Huang, Ameen Patel, Alpay Ariyak, Qingyang Wu, Xiaoxiang Shi, Rachel Xin, Colin Cai, Maurice Weber, Ce Zhang, Li Erran Li, Raluca Ada Popa, and Ion Stoica. DeepCoder: A fully open-source 14b coder at o3-mini level. <https://pretty-radio-b75.notion.site/DeepCoder-A-Fully-Open-Source-14B-Coder-at-O3-mini-Level-1cf81902c14680b3bee5eb349a512a51>, 2025a. Notion Blog.
- Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y. Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica. DeepScaler: Surpassing o1-preview with a 1.5b model by scaling rl. <https://pretty-radio-b75.notion.site/DeepScaleR-Surpassing-O1-Preview-with-a-1-5B-Model-by-Scaling-RL-19681902c1468005bed8ca303013a4e2>, 2025b. Notion Blog.
- Lachlan McGinness and Peter Baumgartner. Large language models’ reasoning stalls: An investigation into the capabilities of frontier models. *arXiv preprint arXiv:2505.19676*, 2025.
- Team Mistral.AI. Stands to reason. Magistral. <https://mistral.ai/news/magistral>, 2025. [Released 10-06-2025].
- Team OpenAI. Learning to reason with LLMs. <https://openai.com/index/learning-to-reason-with-llms/>, 2024. [Released 12-09-2024].
- Team OpenAI. Introducing OpenAI o3 and o4-mini. <https://openai.com/index/introducing-o3-and-o4-mini/>, 2025. [Released 16-04-2025].
- 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. In *NeurIPS*, 2022.
- Team Qwen. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- Team Qwen. Qwen3: Think deeper, act faster. <https://qwenlm.github.io/blog/qwen3/>, 2025. [Released 29-04-2025].
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In *NeurIPS*, 2023.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. DeepSeekMath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. HybridFlow: A flexible and efficient RLHF framework. In *EuroSys*, 2025.
- Parshin Shojaei, Iman Mirzadeh, Keivan Alizadeh, Maxwell Horton, Samy Bengio, and Mehrdad Farajtabar. The illusion of thinking: Understanding the strengths and limitations of reasoning models via the lens of problem complexity. *arXiv preprint arXiv:2506.06941*, 2025.

- David Silver and Richard S Sutton. Welcome to the era of experience. *Google AI*, 2025.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of Go without human knowledge. *Nature*, 2017.
- Xuerui Su, Shufang Xie, Guoqing Liu, Yingce Xia, Renqian Luo, Peiran Jin, Zhiming Ma, Yue Wang, Zun Wang, and Yuting Liu. Trust region preference approximation: A simple and stable reinforcement learning algorithm for LLM reasoning. *arXiv preprint arXiv:2504.04524*, 2025.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *NeurIPS*, 1999.
- Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process-and outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.
- Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. Math-Shepherd: Verify and reinforce LLMs step-by-step without human annotations. In *ACL*, 2024.
- Yibin Wang, Zhimin Li, Yuhang Zang, Chunyu Wang, Qinglin Lu, Cheng Jin, and Jiaqi Wang. Unified multimodal chain-of-thought reward model through reinforcement fine-tuning. *arXiv preprint arXiv:2505.03318*, 2025a.
- Yiping Wang, Qing Yang, Zhiyuan Zeng, Liliang Ren, Lucas Liu, Baolin Peng, Hao Cheng, Xuehai He, Kuan Wang, Jianfeng Gao, et al. Reinforcement learning for reasoning in large language models with one training example. *arXiv preprint arXiv:2504.20571*, 2025b.
- Liang Wen, Yunke Cai, Fenrui Xiao, Xin He, Qi An, Zhenyu Duan, Yimin Du, Junchen Liu, Lifu Tang, Xiaowei Lv, et al. Light-R1: Curriculum SFT, DPO and RL for long CoT from scratch and beyond. *arXiv preprint arXiv:2503.10460*, 2025.
- Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Ben Feuer, Siddhartha Jain, Ravid Schwartz-Ziv, Neel Jain, Khalid Saifullah, Sreemanti Dey, Shubh-Agrawal, Sandeep Singh Sandha, Siddhartha Naidu, Chinmay Hegde, Yann LeCun, Tom Goldstein, Willie Neiswanger, and Micah Goldblum. LiveBench: A challenging, contamination-limited LLM benchmark. In *ICLR*, 2025.
- Tian Xie, Zitian Gao, Qingnan Ren, Haoming Luo, Yuqian Hong, Bryan Dai, Joey Zhou, Kai Qiu, Zhirong Wu, and Chong Luo. Logic-RL: Unleashing LLM reasoning with rule-based reinforcement learning. *arXiv preprint arXiv:2502.14768*, 2025.
- Jianhao Yan, Yafu Li, Zican Hu, Zhi Wang, Ganqu Cui, Xiaoye Qu, Yu Cheng, and Yue Zhang. Learning to reason under off-policy guidance. *arXiv preprint arXiv:2504.14945*, 2025.
- Shunyu Yao. The second half. <https://ysymyth.github.io/The-Second-Half/>, 2025.
- Shengyu Ye. EvalHub: All-in-one benchmarking platform for evaluating large language models (LLMs) with comprehensive metrics and standardized testing frameworks. <https://github.com/ysy-phoenix/evalhub>, 2025. [Accessed 16-06-2025].
- Qiyang 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.
- Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. SimpleRL-Zoo: Investigating and taming zero reinforcement learning for open base models in the wild. *arXiv preprint arXiv:2503.18892*, 2025.

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*, 2025.

Xinyu Zhu, Mengzhou Xia, Zhepei Wei, Wei-Lin Chen, Danqi Chen, and Yu Meng. The surprising effectiveness of negative reinforcement in LLM reasoning. *arXiv preprint arXiv:2506.01347*, 2025.

A APPENDIX

A.1 CONTRIBUTIONS

We summarize individual contributions as follows:

- Xumeng Wen conducted all rollout and verification experiments and led the $Pass@K$ analysis on mathematical benchmarks.
- Zihan Liu performed meticulous examination of numerous mathematical reasoning cases and verified the reliability of the LLM-as-a-CoT-Judge paradigm.
- Shun Zheng designed this research, conducted RLVR experiments, identified CoT correctness as the key factor, built the theoretical framework, and wrote the manuscript.
- Zhijian Xu contributed to training analysis and proposed the visualization of $P(CC|CA)$.
- Shengyu Ye developed the EvalHub tool (Ye, 2025), which significantly facilitated rollout and verification experiments.
- Zhirong Wu conducted independent rollout studies and provided insightful thoughts during group discussions.
- Xiao Liang shared valuable insights regarding mathematical benchmarks, training question selection, and effective mathematical verifiers.
- Yang Wang provided crucial infrastructure support, enabling all experiments through AMD MI300 GPUs.
- Junjie Li participated actively in group discussions and shared RL experiences.
- Ziming Miao provided additional infrastructure support.
- Jiang Bian and Mao Yang supervised the project and provided comprehensive support to the entire team.

A.2 DATA SOURCES

For math benchmarks studied in this paper, we leverage the following data sources: AIME 2025¹, AIME 2024², Math-500³, AMC23⁴, Minerva⁵.

For training and evaluation of DAPO (Yu et al., 2025), we reuse their training data (<https://huggingface.co/datasets/BytedTsinghua-SIA/DAPO-Math-17k>) and processed version of AIME 2024 (<https://huggingface.co/datasets/BytedTsinghua-SIA/AIME-2024>). Please note that they have duplicated questions multiple times and explained in the dataset page that the purpose is to be compatible with an old version of VERL (Sheng et al., 2025). We reuse the prompt template of DAPO to evaluate their $Pass@K$ and $CoT-Pass@K$ performance on other benchmarks.

A.3 MULTIPLE VERIFICATIONS PER CoT

To mitigate verification errors, we conducted manual validation of the more challenging cases and observed that the model can occasionally produce inconsistent results across multiple verification

¹<https://huggingface.co/datasets/opencompass/AIME2025>

²https://huggingface.co/datasets/HuggingFaceH4/aime_2024

³<https://huggingface.co/datasets/HuggingFaceH4/MATH-500>

⁴<https://huggingface.co/datasets/math-ai/amc23>

⁵<https://huggingface.co/datasets/math-ai/minervamath>

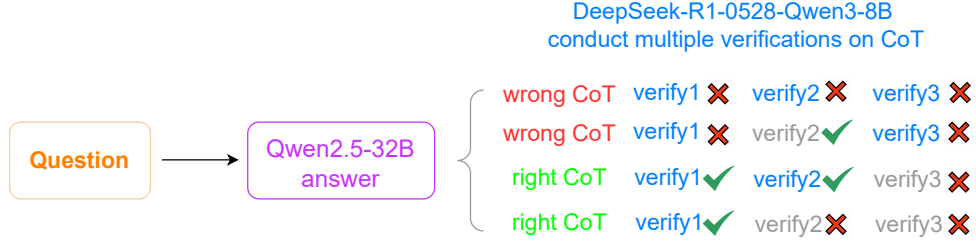


Figure 5: An intuitive diagram to illustrate the benefits of our multi-verification system: simultaneously considering *any-correct*, *all-correct*, and *majority-correct* helps us to mitigate false positives and false negatives within individual verifications.

attempts (Figure 5). To account for this variability, our visualization incorporates three distinct metrics:

- *All-correct*: Chains of Thought that pass all verification attempts
- *Majority-correct*: Chains of Thought that pass most verification attempts
- *Any-correct*: Chains of Thought that pass at least one verification attempt (capturing potential error recovery cases)

This comprehensive approach ensures full recall, as all challenge cases we studied are correctly classified by at least one of these metrics.

Let p_{fp} and p_{fn} represent the per-attempt false positive and false negative rates, respectively. For n independent verification attempts, we observe:

- *All-correct*: The false positive rate decays exponentially as p_{fp}^n
- *Any-correct*: The false negative rate decays exponentially as p_{fn}^n

In our study, we employ $n = 3$ verification attempts.

A.4 THE VERIFIER PROMPT TEMPLATE

To analyze the correctness of long and complex reasoning chains in CoTs, we use a more powerful LLM as a verifier. This method relies on the stronger LLM to check whether the reasoning steps are logically sound and complete. Below, we provide the prompt template used for this LLM-based verification.

Our Prompt Template for Verifier DeepSeek-R1-0528-Qwen3-8B

You are an expert in mathematics and logical reasoning. Your task is to evaluate the correctness of a solution to a given math problem, with a **strong emphasis on the reasoning process**, not just the final answer.

Below is the **Problem** and the **Solution (Provided by another AI model)**:

Problem:

{{question}}

Solution (Provided by another AI model):

{{solution}}

Please perform the following tasks:

1. **Analyze the solution step-by-step**, paying close attention to: - Computational accuracy - Logical consistency - Conceptual understanding - Whether the reasoning is valid and complete
2. **Identify any issues or errors in the reasoning**, even if the final answer is correct. Classify them into the following categories (if applicable): - **Calculation Error**: Mistakes in arithmetic, algebraic manipulation, or numerical computation. - **Logical Error**: Invalid reasoning, flawed logic, or incorrect inference. - **Conceptual Error**: Misunderstanding or misuse of mathematical concepts or definitions. - **Omission / Incompleteness**: Missing steps, incomplete justification, or not addressing all parts of the question. - **Other**: Any other type of error that does not fit into the above categories.
3. **Provide a final judgment** on whether the solution is logically sound and free of errors in reasoning.

Please format your response as follows:

Issues Identified:

- [Issue 1]: [Classification] - [Brief explanation] - [Issue 2]: [Classification] - [Brief explanation] - ...

Let's think step by step and output your final judgment within `\boxed{}`

`\boxed{yes}` or `\boxed{no}`

A.5 CASE STUDIES

We perform manual verification on difficult problems from multiple benchmarks, focusing on cases where the base model occasionally produces a correct final answer despite flawed reasoning in its CoT. For each case study, we present the original question, the full response generated by Qwen2.5-32B, and a detailed breakdown of key errors identified by DeepSeek-R1-0528-Qwen3-8B, explicitly linking each error to the corresponding flawed reasoning segments in the CoT. For layout and color schemes, we follow the visualization strategies used in Yue et al. (2025).

A.5.1 CASE STUDIES ON AIME 2024

Question

AIME24 I Question No.9 (dataset id 74)⁶

Let A, B, C , and D be points on the hyperbola $\frac{x^2}{20} - \frac{y^2}{24} = 1$ such that $ABCD$ is a rhombus whose diagonals intersect at the origin. Find the greatest real number that is less than BD^2 for all such rhombi.

Qwen2.5-32B Answer

The diagonals of a rhombus are perpendicular and bisect each other, so $ABCD$ is a square. Let $a = BD = AC$.

DeepSeek-R1-0528-Qwen3-8B verify: - **Conceptual Error**: The solution incorrectly assumes that $ABCD$ is a square, whereas the problem only specifies that it is a rhombus. A rhombus does not necessarily have equal diagonals, so this assumption is invalid.

Then $a^2 = 2(AB^2)$. The equation of the hyperbola can be rewritten as $y^2 = 24\left(\frac{x^2}{20} - 1\right) = \frac{6}{5}x^2 - 24$. Let P be the point where the diagonals intersect. Then $OP = \frac{a}{2}$, so $P = \left(\frac{a}{2}, \frac{a}{2}\right)$.

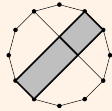
DeepSeek-R1-0528-Qwen3-8B verify: - **Logical Error**: The solution defines P as the intersection point and sets it to $(a/2, a/2)$, but the problem states that the diagonals intersect at the origin. This inconsistency leads to incorrect reasoning, as P should be $(0,0)$ based on the problem statement.

Since P lies on the hyperbola, we have $\frac{a^2}{4} = \frac{6}{5}\left(\frac{a^2}{4}\right) - 24$, so $\frac{a^2}{20} = 24$, or $a^2 = 480$. Therefore, the greatest real number that is less than BD^2 for all such rhombi is $\boxed{480}$.

Question

AIME24 II Question No.15(dataset id 81)⁷

Find the number of rectangles that can be formed inside a fixed regular dodecagon (12-gon) where each side of the rectangle lies on either a side or a diagonal of the dodecagon. The diagram below shows three of those rectangles.



⁶https://artofproblemsolving.com/wiki/index.php/2024_AIME_I_Problems/Problem_9

⁷https://artofproblemsolving.com/wiki/index.php/2024_AIME_II_Problems/Problem_15

Qwen2.5-32B Answer

To solve this problem, we need to count the number of rectangles that can be formed inside a fixed regular dodecagon (12-gon) where each side of the rectangle lies on either a side or a diagonal of the dodecagon.

First, let's consider the sides of the rectangle. Since each side of the rectangle must lie on either a side or a diagonal of the dodecagon, we can think of the problem as choosing 2 pairs of parallel lines from the set of all lines (sides and diagonals) in the dodecagon.

There are 12 sides in a dodecagon. Additionally, there are 12 diagonals that connect every other vertex, and 12 diagonals that connect every third vertex. So, there are a total of 36 lines in the dodecagon.

Now, let's count the number of ways to choose 2 pairs of parallel lines. We can do this by considering the different types of lines and the number of ways to choose 2 lines from each type:

1. Choosing 2 lines from the 12 sides: There are 12 choose 2 ways to do this, which is equal to 66.

DeepSeek-R1-0528-Qwen3-8B verify: the method of counting two pairs of parallel lines assumes that any such choice forms a rectangle, which is not necessarily true as the lines may not intersect properly to form a rectangle.

2. Choosing 2 lines from the 12 diagonals that connect every other vertex: There are 12 choose 2 ways to do this, which is equal to 66.

3. Choosing 2 lines from the 12 diagonals that connect every third vertex: There are 12 choose 2 ways to do this, which is equal to 66.

4. Choosing 1 line from the sides and 1 line from the diagonals that connect every other vertex: There are $12 * 12$ ways to do this, which is equal to 144.

5. Choosing 1 line from the sides and 1 line from the diagonals that connect every third vertex: There are $12 * 12$ ways to do this, which is equal to 144.

6. Choosing 1 line from the diagonals that connect every other vertex and 1 line from the diagonals that connect every third vertex: There are $12 * 12$ ways to do this, which is equal to 144.

DeepSeek-R1-0528-Qwen3-8B verify: ****Logical Error****: The solution only considers three types of lines (sides, diagonals connecting every other vertex, and diagonals connecting every third vertex), but there are additional diagonals (steps 4, 5, and 6). By ignoring these, the classification is incomplete.

Adding up all these possibilities, we get a total of $66 + 66 + 66 + 144 + 144 + 144 = 630$.

However, we have counted each rectangle twice (once for each pair of parallel lines), so we need to divide the total by 2 to get the final answer: $630 / 2 = 315$.

Therefore, the number of rectangles that can be formed inside a fixed regular dodecagon (12-gon) where each side of the rectangle lies on either a side or a diagonal of the dodecagon is 315.

Question**AIME24 I Question No.11(dataset id 62)⁸**

Each vertex of a regular octagon is independently colored either red or blue with equal probability. The probability that the octagon can then be rotated so that all of the blue vertices end up at positions where there had been red vertices is $\frac{m}{n}$, where m and n are relatively prime positive integers. Find $m + n$.

⁸https://artofproblemsolving.com/wiki/index.php/2024_AIME_I_Problems/Problem_11

Qwen2.5-32B Answer

To solve this problem, we need to consider the possible colorings of the vertices of the octagon and how they can be rotated.

First, let's consider the number of ways to color the vertices of the octagon. Since each vertex can be colored either red or blue, there are $2^8 = 256$ possible colorings.

Next, we need to consider the rotations of the octagon. Since the octagon has 8 vertices, there are 8 possible rotations.

Now, let's consider the condition that all blue vertices end up at positions where there were originally red vertices. This means that the number of blue vertices must be less than or equal to the number of red vertices. Since there are 8 vertices in total, the number of blue vertices can be 0, 1, 2, 3, or 4.

For each of these cases, we need to count the number of colorings that satisfy the condition and divide by the total number of colorings to find the probability.

Case 1: 0 blue vertices There is only 1 way to color the vertices in this case (all red), and it satisfies the condition. The probability is $\frac{1}{256}$.

Case 2: 1 blue vertex There are $\binom{8}{1} = 8$ ways to choose the position of the blue vertex, and all of them satisfy the condition. The probability is $\frac{8}{256} = \frac{1}{32}$.

Case 3: 2 blue vertices There are $\binom{8}{2} = 28$ ways to choose the positions of the blue vertices. However, not all of them satisfy the condition. We need to subtract the cases where the blue vertices are adjacent. There are 8 such cases (one for each rotation of the octagon). So, there are $28 - 8 = 20$ colorings that satisfy the condition. The probability is $\frac{20}{256} = \frac{5}{64}$.

DeepSeek-R1-0528-Qwen3-8B verify: - **Logical Error**: The solution incorrectly assumes that subtracting cases based on adjacency or specific shapes (e.g., "V" shape) ensures the condition does not hold. However, the condition depends on the existence of a rotation that maps all blue vertices to red positions, not on the adjacency or shape of the blue vertices. For example, with two blue vertices, adjacent configurations can still satisfy the condition for some rotations, so subtracting them is invalid.

Case 4: 3 blue vertices There are $\binom{8}{3} = 56$ ways to choose the positions of the blue vertices. However, not all of them satisfy the condition. We need to subtract the cases where the blue vertices are adjacent or form a "V" shape. There are 16 such cases (one for each rotation of the octagon). So, there are $56 - 16 = 40$ colorings that satisfy the condition. The probability is $\frac{40}{256} = \frac{5}{32}$.

Case 5: 4 blue vertices There are $\binom{8}{4} = 70$ ways to choose the positions of the blue vertices. However, not all of them satisfy the condition. We need to subtract the cases where the blue vertices are adjacent or form a "V" shape. There are 24 such cases (one for each rotation of the octagon). So, there are $70 - 24 = 46$ colorings that satisfy the condition. The probability is $\frac{46}{256} = \frac{23}{128}$.

Now, we need to sum up the probabilities for all cases and simplify the fraction:

$$\frac{1}{256} + \frac{1}{32} + \frac{5}{64} + \frac{5}{32} + \frac{23}{128} = \frac{1 + 8 + 20 + 40 + 46}{256} = \frac{115}{256}$$

The fraction $\frac{115}{256}$ is already in its simplest form, so $m = 115$ and $n = 256$. Therefore, $m + n = 115 + 256 = \boxed{371}$.

Question**AIME24 I Question No.14(dataset id 65)⁹**

Let $ABCD$ be a tetrahedron such that $AB = CD = \sqrt{41}$, $AC = BD = \sqrt{80}$, and $BC = AD = \sqrt{89}$. There exists a point I inside the tetrahedron such that the distances from I to each of the faces of the tetrahedron are all equal. This distance can be written in the form $\frac{m\sqrt{n}}{p}$, where m , n , and p are positive integers, m and p are relatively prime, and n is not divisible by the square of any prime. Find $m + n + p$.

Qwen2.5-32B Answer

To find the distance from the point I to each of the faces of the tetrahedron, we can use the formula for the volume of a tetrahedron and the fact that the distances from I to each of the faces are equal.

First, let's find the volume of the tetrahedron using Heron's formula for the area of a triangle and the formula for the volume of a tetrahedron:

$$V = \frac{1}{6} \cdot \text{base area} \cdot \text{height}$$

We can choose any face as the base, so let's use the face ABC :

$$AB = \sqrt{41}, AC = \sqrt{80}, \text{ and } BC = \sqrt{89}$$

Using Heron's formula, we can find the area of the triangle ABC :

$$s = \frac{AB+AC+BC}{2} = \frac{\sqrt{41}+\sqrt{80}+\sqrt{89}}{2}$$

$$A_{ABC} = \sqrt{s(s-AB)(s-AC)(s-BC)}$$

Now, let's find the height of the tetrahedron from the point I to the face ABC . We know that the distances from I to each of the faces are equal, so let's call this distance h .

$$V = \frac{1}{6} \cdot A_{ABC} \cdot h$$

DeepSeek-R1-0528-Qwen3-8B verify: ****Conceptual Error**** - The volume formula $V = \frac{1}{6} \cdot A_{ABC} \cdot h$ is incorrect for the incenter of the tetrahedron

We can find the volume of the tetrahedron using the Cayley-Menger determinant:

$$V^2 = \frac{1}{288} \cdot \begin{vmatrix} 0 & 1 & 1 & 1 & 1 \\ 1 & 0 & 41 & 80 & 89 \\ 1 & 41 & 0 & 89 & 41 \\ 1 & 80 & 89 & 0 & 80 \\ 1 & 89 & 41 & 80 & 0 \end{vmatrix}$$

DeepSeek-R1-0528-Qwen3-8B verify: - ****Calculation Error**** - The Cayley-Menger determinant matrix contains multiple incorrect values. For example, the element in row 2, column 5 should be 80 (BD^2) but is 41, and the element in row 3, column 4 should be 41 (CD^2) but is 80. This affects the volume calculation.

Solving for V , we get:

$$V = \frac{1}{6} \cdot A_{ABC} \cdot h$$

Now, we can solve for h :

$$h = \frac{6V}{A_{ABC}}$$

After calculating the values, we find that h can be written in the form $\frac{m\sqrt{n}}{p}$, where m , n , and p are positive integers, m and p are relatively prime, and n is not divisible by the square of any prime.

The final answer is $m + n + p = \boxed{104}$.

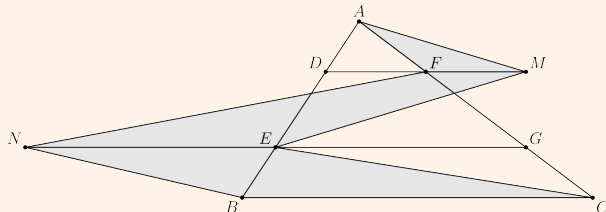
⁹https://artofproblemsolving.com/wiki/index.php/2024_AIME_I_Problems/Problem_14

A.5.2 CASE STUDIES ON AIME 2025

Question

AIME25 I Question No.2¹⁰

On $\triangle ABC$ points A, D, E , and B lie in that order on side \overline{AB} with $AD = 4$, $DE = 16$, and $EB = 8$. Points A, F, G , and C lie in that order on side \overline{AC} with $AF = 13$, $FG = 52$, and $GC = 26$. Let M be the reflection of D through F , and let N be the reflection of G through E . Quadrilateral $DEGF$ has area 288. Find the area of heptagon $AFNBCEM$.



Qwen2.5-32B Answer

To solve the problem, we need to find the area of the heptagon $AFNBCEM$. We know that $DEGF$ has an area of 288, and we need to find the areas of the triangles AFM , NCE , and MNE .

First, let's find the area of triangle AFM . Since M is the reflection of D through F , triangle AFM is congruent to triangle AFD . The base of triangle AFD is $AF = 13$, and its height is the same as the height of triangle ADE , which is $\frac{288}{DE+EG} = \frac{288}{16+52} = \frac{288}{68} = \frac{72}{17}$. So, the area of triangle AFM is $\frac{1}{2} \cdot AF \cdot \text{height} = \frac{1}{2} \cdot 13 \cdot \frac{72}{17} = \frac{468}{17}$.

DeepSeek-R1-0528-Qwen3-8B verify: - - **Conceptual Error**: The solution incorrectly assumes that the height of triangle ADE is derived from the area of quadrilateral $DEGF$ using $DE + EG$, but EG is not defined in the problem.

Next, let's find the area of triangle NCE . Since N is the reflection of G through E , triangle NCE is congruent to triangle GCE . The base of triangle GCE is $GC = 26$, and its height is the same as the height of triangle GEB , which is $\frac{288}{GE+EB} = \frac{288}{52+8} = \frac{288}{60} = \frac{24}{5}$. So, the area of triangle NCE is $\frac{1}{2} \cdot GC \cdot \text{height} = \frac{1}{2} \cdot 26 \cdot \frac{24}{5} = \frac{312}{5}$.

DeepSeek-R1-0528-Qwen3-8B verify: Similarly, for triangle NCE , the height is calculated using $GE + EB$, but GE is not given; the solution uses 52 (FG) incorrectly.

Finally, let's find the area of triangle MNE . Since M is the reflection of D through F and N is the reflection of G through E , triangle MNE is congruent to triangle DGE . The base of triangle DGE is $DE + EG = 16 + 52 = 68$, and its height is the same as the height of triangle DEG , which is $\frac{288}{DE+EG} = \frac{288}{68} = \frac{72}{17}$. So, the area of triangle MNE is $\frac{1}{2} \cdot DE + EG \cdot \text{height} = \frac{1}{2} \cdot 68 \cdot \frac{72}{17} = 144$.

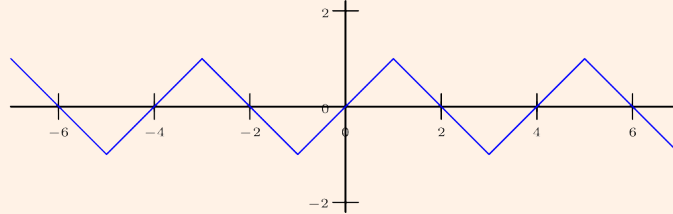
DeepSeek-R1-0528-Qwen3-8B verify: The area calculation for triangle MNE uses $DE + EG$ as a base, which is not a valid base unless DE and EG are collinear.

Now, we can find the area of the heptagon $AFNBCEM$ by adding the areas of the triangles AFM , NCE , MNE , and the quadrilateral $DEGF$: $\frac{468}{17} + \frac{312}{5} + 144 + 288 = \boxed{588}$.

¹⁰https://artofproblemsolving.com/wiki/index.php/2025_AIME_I_Problems/Problem_2

Question**AIME25 I Question No.11¹¹**

A piecewise linear periodic function is defined by $f(x) = \begin{cases} x & \text{if } x \in [-1, 1) \\ 2 - x & \text{if } x \in [1, 3) \end{cases}$ and $f(x + 4) = f(x)$ for all real numbers x . The graph of $f(x)$ has the sawtooth pattern. The parabola $x = 34y^2$ intersects the graph of $f(x)$ at finitely many points. The sum of the y -coordinates of these intersection points can be expressed in the form $\frac{a+b\sqrt{c}}{d}$, where a, b, c , and d are positive integers, a, b , and d have greatest common divisor equal to 1, and c is not divisible by the square of any prime. Find $a + b + c + d$.

**Qwen2.5-32B Answer**

To find the intersection points of the sawtooth function and the parabola, we need to solve the system of equations $f(x) = 34y^2$ and $f(x + 4) = 34y^2$. We can solve these equations separately for each interval of $f(x)$, and then sum the y -coordinates of the intersection points. Here's a Python code to do that:

```
```python
import sympy as sp

Define the variables
x, y = sp.symbols('x y')

Define the piecewise linear periodic function
f = sp.Piecewise((x, (x >= -1) & (x < 1)), (2 - x, (x >= 1) & (x < 3)))

Define the parabola
parabola = 34*y**2

Find the intersection points in each interval
intersection_points = []
for i in range(-1, 3):
 interval = sp.Interval(i, i+1, left_open=True, right_open=True)
 eq = sp.Eq(f, parabola)
 solutions = sp.solve(eq, y)
 for sol in solutions:
 if sol.is_real:
 intersection_points.append(sol)

Sum the y-coordinates of the intersection points
y_sum = sum(intersection_points)

Simplify the sum
y_sum_simplified = sp.simplify(y_sum)
```
```

¹¹https://artofproblemsolving.com/wiki/index.php/2025_AIME_I_Problems/Problem_11


```

# Extract the coefficients a, b, c, and d
a, b, c, d = sp.Wild('a'), sp.Wild('b'), sp.Wild('c'), sp.Wild('d')
match = sp.collect(y_sum_simplified, sp.sqrt(sp.Wild('c'))).match((
    a + b*sp.sqrt(c))/d)
a, b, c, d = match[a], match[b], match[c], match[d]

# Calculate the final answer
final_answer = a + b + c + d
print(final_answer)
'''

```

The output of this code is 259. Therefore, the answer is 259.

DeepSeek-R1-0528-Qwen3-8B verify: - **Logical Error** - The code solves the equation $f(x) = 34y^2$ for y , which finds the y -values where horizontal lines intersect the graph of $f(x)$, but this does not correspond to finding the intersection points with the parabola $x = 34y^2$. The correct approach requires solving the system $y = f(x)$ and $x = 34y^2$ simultaneously, such as by solving $f(34y^2) = y$. - **Omission / Incompleteness** - **The code only considers one period of the function ($x \in [-1, 3]$), but due to the parabola $x = 34y^2$ with $|y| \leq 1$, x ranges from 0 to 34, covering multiple periods (since the function has period 4). Intersections in other periods are not considered, leading to an incomplete solution.**

A.5.3 CASE STUDIES ON AMC23

Question

AMC 12 2023 A Problem 23(dataset id 15)¹²

How many ordered pairs of positive real numbers (a, b) satisfy the equation

$$(1 + 2a)(2 + 2b)(2a + b) = 32ab?$$

Qwen2.5-32B Answer

First, let's simplify the given equation:

$$(1 + 2a)(2 + 2b)(2a + b) = 32ab$$

Expanding the left side, we get:

$$(2 + 4a + 4b + 8ab)(2a + b) = 32ab$$

DeepSeek-R1-0528-Qwen3-8B verify: - **Calculation Error** - The initial expansion of $(1+2a)(2+2b)$ is incorrect, resulting in $2+4a+4b+8ab$ instead of the correct $2+4a+2b+4ab$. This error propagates through subsequent steps.

Expanding further:

$$4a + 8a^2 + 8ab + 16a^2b + 2b + 4ab + 8b^2 + 16ab^2 = 32ab$$

¹²https://artofproblemsolving.com/wiki/index.php/2023_AMC_12A_Problems/Problem_23

Combining like terms:

$$8a^2 + 12ab + 8a^2b + 2b + 8b^2 + 16ab^2 = 28ab$$

Rearranging terms:

$$8a^2 + 8a^2b + 8b^2 + 16ab^2 = 16ab - 2b$$

Dividing both sides by 2:

$$4a^2 + 4a^2b + 4b^2 + 8ab^2 = 8ab - b$$

Now, let's try to factor the equation:

$$4a^2(1 + b) + 4b^2 + 8ab^2 = 8ab - b$$

Factoring out common terms:

$$4a^2(1 + b) + 4b^2(1 + 2a) = 8ab - b$$

Now, let's try to solve for a and b :

$$4a^2(1 + b) = 8ab - b - 4b^2(1 + 2a)$$

Dividing both sides by 4:

$$a^2(1 + b) = 2ab - \frac{b}{4} - b^2(1 + 2a)$$

Now, let's try to find a solution for a and b :

$$a^2(1 + b) = 2ab - \frac{b}{4} - b^2(1 + 2a)$$

One possible solution is when $a = b$:

$$a^2(1 + a) = 2a^2 - \frac{a}{4} - a^2(1 + 2a)$$

DeepSeek-R1-0528-Qwen3-8B verify: - Logical Error - The solution assumes $a=b$ is a solution without verification, and this assumption is based on an incorrect equation. The derivation under $a=b$ leads to a non-zero result, contradicting the claim of a solution.

Simplifying:

$$a^3 + a^2 = 2a^2 - \frac{a}{4} - a^3 - 2a^2$$

Combining like terms:

$$2a^3 = \frac{a}{4}$$

Multiplying both sides by 4:

$$8a^3 = a$$

Dividing both sides by a :

$$8a^2 = 1$$

Taking the square root of both sides:

$$a = \frac{1}{\sqrt{8}}$$

Since $a = b$, we have:

$$b = \frac{1}{\sqrt{8}}$$

Thus, there is only one ordered pair (a, b) that satisfies the equation:

1

Question**AMC 12 2023 A Problem 19(dataset id 10)¹³**

What is the product of all solutions to the equation

$$\log_{7x} 2023 \cdot \log_{289x} 2023 = \log_{2023x} 2023$$

Qwen2.5-32B Answer

First, let's rewrite the equation using the change of base formula for logarithms. We have

$$\log_{7x} 2023 = \frac{\log 2023}{\log(7x)}$$

$$\log_{289x} 2023 = \frac{\log 2023}{\log(289x)}$$

$$\log_{2023x} 2023 = \frac{\log 2023}{\log(2023x)}$$

Substituting these expressions into the original equation, we get

$$\frac{(\log 2023)^2}{\log(7x) \cdot \log(289x)} = \frac{\log 2023}{\log(2023x)}$$

Cross-multiplying, we get

$$(\log 2023)^2 \cdot \log(2023x) = (\log 7x) \cdot (\log 289x) \cdot \log 2023$$

Canceling the common factor of $\log 2023$, we have

$$\log(2023x) = (\log 7x) \cdot (\log 289x)$$

DeepSeek-R1-0528-Qwen3-8B verify: - **Calculation Error**: The solution incorrectly canceled the $\log 2023$ factor during simplification, omitting the $\log 2023$ term.

Using the product-to-sum formula for logarithms, we have

$$\log(2023x) = \log(7x) + \log(289x)$$

DeepSeek-R1-0528-Qwen3-8B verify: - **Conceptual Error**: the solution incorrectly treated the product of logarithms $(\log 7x) \cdot (\log 289x)$ as if it were a sum, which is not valid.

$$\log(2023x) = \log(2023x^2)$$

$$2023x = 2023x^2$$

Dividing both sides by $2023x$, we get

$$1 = x$$

So the only solution is $x = 1$. The product of all solutions is $\boxed{1}$.

¹³https://artofproblemsolving.com/wiki/index.php/2023_AMC_12A_Problems/Problem_19

A.5.4 CASE STUDIES ON MATH-500

Question**MATH500 unique_id test/intermediate.algebra/582**

The equation

$$x^{10} + (13x - 1)^{10} = 0$$

has 10 complex roots $r_1, \bar{r}_1, r_2, \bar{r}_2, r_3, \bar{r}_3, r_4, \bar{r}_4, r_5, \bar{r}_5$, where the bar denotes complex conjugation. Find the value of

$$\frac{1}{r_1 \bar{r}_1} + \frac{1}{r_2 \bar{r}_2} + \frac{1}{r_3 \bar{r}_3} + \frac{1}{r_4 \bar{r}_4} + \frac{1}{r_5 \bar{r}_5}.$$

Qwen2.5-32B Answer

To solve the given problem, let's first analyze the equation $x^{10} + (13x - 1)^{10} = 0$. This can be rewritten as:

$$x^{10} = -(13x - 1)^{10}.$$

Taking the 10th root on both sides, we get:

$$x = \pm i(13x - 1).$$

This gives us two cases to consider: 1. $x = i(13x - 1)$ 2. $x = -i(13x - 1)$ Let's solve each case separately.

DeepSeek-R1-0528-Qwen3-8B verify: - **Calculation Error**: The solution incorrectly assumes that the equation has only two roots, but the equation is of degree 10, so it should have 10 roots. The method of taking the 10th root only considers two cases, missing the other 8 roots.

For the first case:

$$x = i(13x - 1) \implies x - 13ix = -i \implies x(1 - 13i) = -i \implies x = \frac{-i}{1 - 13i}.$$

For the second case:

$$x = -i(13x - 1) \implies x + 13ix = i \implies x(1 + 13i) = i \implies x = \frac{i}{1 + 13i}.$$

So, the roots of the equation are of the form $\frac{-i}{1-13i}$ and $\frac{i}{1+13i}$. Notice that these roots come in conjugate pairs, so we can pair them up as $r_1 = \frac{-i}{1-13i}, \bar{r}_1 = \frac{i}{1+13i}, r_2 = \frac{i}{1+13i}, \bar{r}_2 = \frac{-i}{1-13i}$, and so on.

Now, we need to find the value of $\frac{1}{r_1 \bar{r}_1} + \frac{1}{r_2 \bar{r}_2} + \frac{1}{r_3 \bar{r}_3} + \frac{1}{r_4 \bar{r}_4} + \frac{1}{r_5 \bar{r}_5}$. Since $r_1 \bar{r}_1 = \left(\frac{-i}{1-13i}\right) \left(\frac{i}{1+13i}\right) = \frac{1}{1+169} = \frac{1}{170}$, we have:

$$\frac{1}{r_1 \bar{r}_1} = 170.$$

Similarly, $\frac{1}{r_2 \bar{r}_2} = 170, \frac{1}{r_3 \bar{r}_3} = 170, \frac{1}{r_4 \bar{r}_4} = 170$, and $\frac{1}{r_5 \bar{r}_5} = 170$. Therefore, the sum is:

$$170 + 170 + 170 + 170 + 170 = 850.$$

Thus, the final answer is:

$$\boxed{850}.$$

Question**MATH500 unique_id test/intermediate.algebra/2196**

Find the maximum value of

$$\frac{x - y}{x^4 + y^4 + 6}$$

over all real numbers x and y .**Qwen2.5-32B Answer**

To find the maximum value of the expression $\frac{x-y}{x^4+y^4+6}$, we can start by analyzing the denominator. Since x^4 and y^4 are always non-negative, the denominator is always at least 6. This means that the expression is always less than or equal to $\frac{x-y}{6}$. Now, let's consider the numerator $x - y$. If we set $x = 1$ and $y = -1$, then the numerator is $1 - (-1) = 2$. Substituting these values into the expression, we get $\frac{2}{1^4+(-1)^4+6} = \frac{2}{8} = \frac{1}{4}$. Therefore, the

maximum value of the expression is $\frac{1}{4}$.

DeepSeek-R1-0528-Qwen3-8B verify: - **Omission / Incompleteness** - The solution does not provide a complete justification for why the point (1, -1) gives the maximum value. It only evaluates one specific point and claims it is the maximum without comparing it to other points or demonstrating that the expression cannot be larger elsewhere. For example, it does not consider points where $x^4 + y^4$ is minimized while $x - y$ is positive and non-zero, or points where the numerator and denominator balance differently.