TOWARDS UNDERSTANDING METACOGNITION IN LARGE REASONING MODELS

Anonymous authorsPaper under double-blind review

000

001

002003004

010 011

012

013

014

016

018

019

021

025

026

027 028 029

030

032

033

034

037

038

040

041 042

043

044

046

047

048

049

051

052

ABSTRACT

Large Reasoning Models (LRMs) have demonstrated remarkable capabilities on complex tasks. Despite these advances, we identify a fundamental limitation: current LRMs impose fixed cognition patterns, lacking the intrinsic ability to be aware of, or regulate their own reasoning processes. This signifies a critical absence of metacognition—an essential faculty in human intelligence. Building on psychology and cognitive science, we first construct a functional framework for metacognition in LRMs, separating internal informational signals from behavioral abilities. This framework is then applied to a comprehensive investigation on seven state-of-the-art LRMs and reveals a consistent gap: while metacognitive information is present and predictive, it often fails to translate into reliable monitoring or control behaviors. To address this gap, we introduce two distinct paradigms for instilling metacognition in LRMs: (1) an emergent approach that leverages prompting to orchestrate metacognitive functions, such as task assessment, confidence monitoring, and strategy regulation; (2) an intrinsic approach that internalizes these faculties by encoding structured meta-cognitive information directly into the model's parameters through training. Overall, our results indicate that integrating metacognitive reasoning improves task performance and offers a valuable lens for the design of future reasoning models.

1 Introduction

Large Reasoning Models (LRMs) like OpenAI-o1 (Jaech et al., 2024) and Deepseek-R1 (Guo et al., 2025a) have achieved remarkable success in complex domains such as coding and mathematics. At first glance, these models appear to exhibit advanced, reflective behaviors within their long chain-of-thought (CoT) (Wei et al., 2022) reasoning. However, a closer look reveals a potential fragility in their cognitive processes. For instance, when tasked with solving a complex-variable equation under the explicit constraint that "z is a positive real number," an LRM may persist in a fixed reasoning pattern like "the problem likely intended for z to be a complex number," a phenomenon termed *reasoning rigidity* (Jang et al., 2025; Araya, 2025). Furthermore, current LRMs also frequently display illusory self-correction, employing introspective phrases like "Wait, let me double-check..." without any actual adjustment to a flawed reasoning trajectory (Guo et al., 2025a; Wang et al., 2025b). These consistent failures in self-monitoring and adaptation reveal a core limitation of current LRMs.

We argue that this deficit can be productively framed as an absence of metacognition (Ackerman & Thompson, 2017; Norman et al., 2019; Tankelevitch et al., 2024)—the ability to monitor and control one's own cognitive processes. In the theory of human cognition, metacognition is essential for evaluating potential reasoning errors (Yeung & Summerfield, 2012), calibrating uncertainty in decision-making (Qiu et al., 2018), and dynamically adapting strategies based on performance (Cary & Reder, 2002). This comparison to human intelligence raises a pivotal question for AI: do current LRMs possess any analogous capabilities, and if so, are they functionally engaged during inference? In this paper, we present the first systematic investigation into this fundamental question.

To structure this investigation, we introduce a functional framework for LRM metacognition, inspired by foundational models in cognitive science (Efklides & Misailidi, 2010; Ackerman & Thompson, 2017). Our framework decomposes metacognition into two components: **information** and **abilities**. Metacognitive information (Dayan, 2023; Norman et al., 2019), the basis for judgment, which includes both static knowledge (e.g., learned strategies in parameters) and dynamic ex-

perience (e.g., internal computational signals). Metacognitive abilities (Nelson & Dunlosky, 1991; Fiedler et al., 2019), the actions taken upon this information, which include *monitoring* (e.g., assessing task difficulty and confidence) and *control* (e.g., selecting reasoning strategies or decomposing problems). By deconstructing metacognition into these components, our framework provides a principled foundation to systematically probe whether, and in what way, metacognition emerges in contemporary LRMs.

Our investigation begins by empirically grounding the first component of our framework: **metacognitive information**. Focusing on the dynamic aspect of *experience*, we probe open-source LRMs to determine whether internal computational signals correlate with reasoning outcomes (§3). Our analysis yields a striking finding: signals spanning the entire Transformer architecture—from input-layer attributions to final-layer token probabilities—are highly predictive of answer success. Critically, we demonstrate that correct and incorrect reasoning traces generate statistically distinguishable internal signatures, providing the first empirical evidence that a machine-readable basis for metacognitive experience exists within these models.

This informational foundation compels the subsequent question: do state-of-the-art LRMs functionally leverage this information as observable metacognitive abilities (§4). Our evaluation of current leading reasoning models across a series of monitoring and control tasks reveals consistent failures. Specifically, we find that the models systematically misjudge task difficulty, display poorly calibrated confidence, and lack proactive planning and strategic flexibility. This evidence suggests that, while predictive metacognitive information may exist internally, it does not reliably translate into effective monitoring or control, exposing a critical gap in the capabilities of current LRM.

To bridge this information-to-ability gap, we propose two complementary paradigms for enhancing LRM metacognition: **①** Emergent Metacognition, scaffolds this pathway at inference time through prompt-guided role-playing. We assign distinct metacognitive roles—such as 'Planner', 'Solver', and 'Verifier'—to simulate a complete monitoring and control loop. This external scaffolding forces the model to act on its latent experiences, effectively eliciting robust metacognitive behaviors without any parameter updates. **②** Internalized Metacognition, directly enriches model's metacognitive knowledge through fine-tuning. We construct a dataset with explicit metacognitive annotations (e.g., plans, self-corrections) and fine-tune the model using a hybrid learning objective, directly embedding these capabilities into its parameters. Together, these two paradigms provide a comprehensive roadmap toward more introspective and reliable reasoning systems.

In summary, our findings reveal a clear dissociation between internal metacognitive information and externally observable metacognitive ability in current LRMs. This gap illuminates a new frontier for research: desgining systems that not only possesses self-awareness but cna also act upon it, effectively bridging latent experience with adaptive behavior.

2 DEFINING METACOGNITION IN LRMS

Metacognition, first conceptualized in developmental psychology, refers to the capacity to monitor and control one's own cognitive processes (Flavell, 1979). According to the well-established two-level model of Nelson and Narens (Nelson & Dunlosky, 1991), metacognition comprises a *object-level* cognition (the act of thinking, perceiving, or remembering) and *meta-level* cognition (the act of thinking about one's thinking). While Large Reasoning Models (LRMs) do not possess subjective consciousness, their complex, multi-step reasoning processes create the functional necessity for such meta-level oversight. We therefore adopt a **functionalist perspective**: we investigate whether LRMs can exhibit behaviors and leverage internal signals that are functionally equivalent to human metacognition, enabling them to produce more reliable and robust reasoning.

Following established frameworks in cognitive science (Tankelevitch et al., 2024), we structure our functional model of LRM metacognition into two core components: Information and Abilities.

Metacognitive Information serves as the basis for judgment. It comprises (i) static *knowledge*—the latent understanding of tasks, strategies, and its own capabilities implicitly encoded in its parameters—and (ii) dynamic *experience*, which we operationalize as the internal computational signals (e.g., token probabilities) generated during a reasoning trace, serving as a functional analogue to a human's 'feeling of error'.

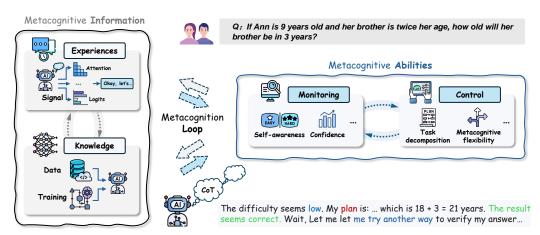


Figure 1: A functional framework for LRM metacognition.

Metacognitive Abilities are the actions taken based on this information. They consist of (i) *monitoring*, the capacity to generate self-assessments about its cognitive state or the task at hand, such as evaluating problem difficulty (§4.1) or estimating its confidence (§4.2); and (ii) *control*, the capacity to strategically alter its reasoning process, such as by performing task decomposition (§4.3) or exhibiting cognitive flexibility when encountering errors (§4.4).

This framework provides a structured lens through which we can systematically investigate the nascent metacognitive capabilities of modern LRMs (i.e., Fig. 1).

3 METACOGNITIVE INFORMATION: KNOWLEDGE AND EXPERIENCES

We begin at the foundation of our proposed framework: **Metacognitive Information**. For an LRM to monitor or control its reasoning, it must first possess information about its own process. This information comprises: (i) *static knowledge*, the vast, latent strategies encoded within the model's parameters, and (ii) *dynamic experience*, the internal information that can directly experience during the reasoning process, such as token probabilities and attention patterns. While static knowledge is the inherent properties of an LRM that are fixed after pre-training, dynamic experience, however, is task-specific, and thereby can contribute to dissecting its correlation with the reasoning correctness. By delving into the internals of an open-source model, we seek to provide the foundational evidence that the information necessary for metacognitive abilities is not only present but also machine-readable, paving the way for the behavioral investigations that follow.

Setup. To test this hypothesis, we utilize three standard benchmarks: GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021), and AIME. We conduct our analysis on Qwen-32B (Yang et al., 2025), a powerful open-source LRM that grants us full access to its internal states. To generate a diverse set of both correct and incorrect reasoning traces for comparison, we use a high temperature T=1.0 for generation, promoting exploration.

Experiment Details. To capture a holistic view of the model's internal state, we consider four types of signals that span the entire processing pipeline of a Transformer block—from input-level importance to output-level confidence. These signals are: (a) Softmax probabilities from the final layer, reflecting output uncertainty; (b) Fully-connected activations and (c) Self-attention scores from the intermediate hidden layers, representing the core of the model's computational state; and (d) Integrated Gradients (IG) attributions at the input layer, indicating perceived input importance. Our analysis is twofold: we first use t-SNE projections for a qualitative visualization of the separability between correct and incorrect samples based on these signals. We then conduct a rigorous quantitative validation by training simple linear classifiers (a 'prober') to predict the final correctness of a trace using only these internal signals as features. See Appendix A.1 for more details.

Results. Qualitatively, we find the internal signal distributions exhibit clear separability for correct (blue) and incorrect (orange) traces (Fig. 2), which provide strong evidence that an LRM's internal signals act as reliable correlates of its reasoning outcomes. Quantitatively, this separation is further confirmed in Tab. 1. Trained solely on these internal signals, the probers can predict the final cor-

Table 1: AUC of the trained linear prober in justifying the reasoning trace correctness.

source	GSM8K	MATH500	AIME2024	AIME2025
Softmax probabilities	0.81	0.68	0.50	0.43
Fully-connected activations	0.79	0.73	0.53	0.50
Self-attention scores	0.71	0.73	0.57	0.53
Integrated Gradient	0.61	0.55	0.40	0.37

rectness of a reasoning trace with an AUC score significantly above chance. These results establish that the signals are both distinct and highly predictive, thereby validating our initial hypothesis.

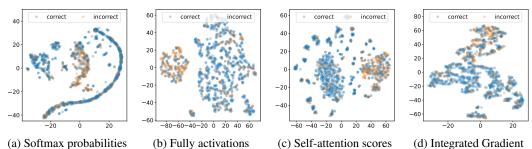


Figure 2: The t-SNE of the internal signals for the first tokens. We capture the activations and attention from the last layer. The distributions are different between the correct and incorrect traces.

4 MEASURING METACOGNITIVE ABILITIES IN LRMS

The evidence of internal metacognitive information in §3 motivates our central question: *can this latent information manifest as observable, functional abilities?* To answer this, we shift our focus from internal correlates to external actions. We thus introduce a new benchmark to systematically measure these abilities across a range of state-of-the-art LRMs from leading developers.

MetaEval. We investigate to what extent LRMs can explicitly monitor and control their own reasoning processes. It is structured around our two-part metacognitive framework, assessing: (1) Metacognitive Monitoring (§4.1, 4.2), probed via self-awareness and confidence adjustment tasks; and (2) Metacognitive Control (§4.3, 4.4), probed via task decomposition and metacognitive flexibility challenges. Unlike existing evaluation suites that focus almost exclusively on object-level task accuracy, MetaEval provides the first targeted evaluation of these crucial, second-order reasoning skills that underpin reliable intelligence.

Models. We examine seven state-of-the-art LRMs to assess the prevalence of these abilities across the AI landscape: Gemini-2.5-Pro, GPT-OSS-120B (Agarwal et al., 2025), Seed-1.5-VL-Pro (Guo et al., 2025b), Doubao-1.5-Pro (Seed et al., 2025), Kimi-K2 (Team et al., 2025), Deepseek-R1 (Guo et al., 2025a), Qwen3-8B/32B (Yang et al., 2025). We sample using temperature T=0.6 for both reasoning and knowledge QA tasks.

4.1 METACOGNITIVE MONITORING: SELF AWARENESS

First, we measure a key aspect of metacognitive monitoring: *self-awareness*. Intuitively, an expert AI reasoner should assess a problem's intrinsic difficulty before committing to a specific solution. This initial assessment allows for the allocation of appropriate cognitive resources and the selection of a suitable strategy. We thus operationalize self-awareness as the model's ability to accurately classify the difficulty of mathematical problems when explicitly prompted to do so.

Experiment Details. To investigate this phenomenon, we design a multi-class classification task from the DEEPMATH103K (He et al., 2025) dataset. We categorize the problems into three primary difficulty levels: Easy (rating < 3.5), Medium ($3.5 \le \text{rating} \le 6.5$), and Hard (rating > 6.5). To elicit difficulty assessments, we prompt each model with "... your task is to assess the difficulty of a math problem based on the provided rubric and examples." Our primary metric is difficulty assessment accuracy, defined as the percentage of problems where a model's predicted category cor-

rectly matches at least one of its ground-truth labels. In addition to the overall evaluation, we report per-category accuracy to analyze model performance on each difficulty level independently. See Appendix A.1.1 for more details.

Results. As shown in Fig. 3, we find the SOTA LRMs exhibit a nascent but highly variable capacity for self-awareness. Top models like <code>Gemini-2.5-Pro</code> achieve an awareness of over 70%, demonstrating a significant ability to distinguish problem complexities. In contrast, other models lag considerably, indicating this is a challenging metacognitive skill. Analysis of the confusion matrices reveals a common failure mode across all models: a strong tendency to misclassify <code>Medium</code> problems as <code>Hard</code>. This suggests a generally conservative or risk-averse assessment strategy, where models are more likely to overestimate than underestimate difficulty when faced with uncertainty. This initial finding demonstrates that while self-awareness can be elicited, its reliability is far from guaranteed.

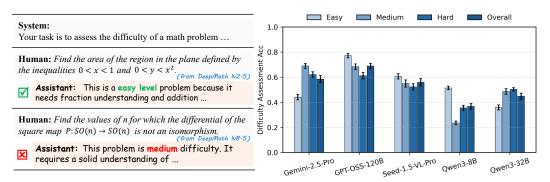


Figure 3: Assessment of LRMs' self-awareness on task complexity. The results show that the popular state-of-the-art LRMs exhibit a lack of capacity in perceiving task difficulty.

4.2 METACOGNITIVE MONITORING: CONFIDENCE AND ITS ADJUSTMENT

We next investigate whether an LRM exhibits metacognitive monitoring by tracking and adjusting its internal confidence during the reasoning process. This ability is crucial, as it allows the system to distinguish correct from incorrect reasoning and signal when its output is untrustworthy.

Experiment Details. We conduct our analysis on challenging benchmarks requiring long-form reasoning, including subsets of DEEPMATH103K, AIME, and GPQA datasets. To quantify the model's internal confidence, we adopt a standard logits-based metric from prior work (Fu et al., 2025), token confidence C_t as the negative average log-probability of the top-k tokens at position t. These token-level scores are then aggregated to produce a trace-level metric, termed of average trace confidence, for each complete solution.

To capture confidence dynamics, we calculate the average trace confidence focusing on the start, middle, and final portions (e.g., 2048 tokens). We then consider four dynamic confidence patterns: consistently high/low (all confidence above/below a high threshold), increasing/decreasing (confidence rises/falls significantly from start to end). We then quantify misalignment between confidence trends and actual correctness: for example, if confidence rises steadily but the final answer is wrong, this indicates poor metacognitive adjustment. The frequency of these events serves as our primary metric for poor confidence adjustment.

Results. We find that LRMs frequently display confidence trajectories that do not correspond with answer correctness (Fig. 4), suggesting weak metacognitive calibration. Notably, confidence often increases even in incorrect traces—indicating that the model becomes more certain as it reasons incorrectly. This gap necessitates the better regulated confidence in alignment with reasoning quality.

4.3 METACOGNITIVE CONTROL: TASK DECOMPOSITION

We now turn to control behavior and examine LRM's task decomposition ability to engage in internal planning prior to reasoning, aiming to determine whether LRMs possess intrinsic planning ability.

Experiment Details. We evaluate on subsets of DEEPMATH, AIME, and GPQA. For each question, we first compute a baseline accuracy using standard direct reasoning prompts. We then

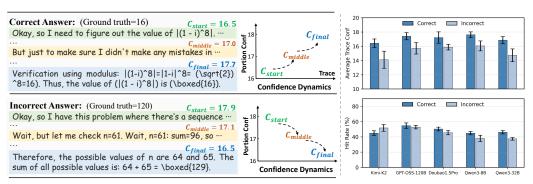


Figure 4: Statistics of internal confidence and adjustment. The average trace confidence of correct and incorrect samples exhibits a significant difference, showing potential as a clear discriminatory signal. Meanwhile, the dynamic confidence patterns suggest weak metacognitive calibration.

introduce task decomposition interventions designed to elicit a plan-before-solve strategy. In the single-turn condition, the prompt instructs the LRM: "Your task is to first break down the problem into a clear, step-by-step plan. Then, execute your plan, reasoning step by step." In the multi-turn condition, the LRM is first asked: "Your ONLY task is to create a high-level, step-by-step plan to solve the following problem." After generating the plan, the plan and original question are concatenated and fed back into the model to complete the reasoning. We compare the final accuracy across these settings to assess whether explicit decomposition enhances reasoning performance, thus revealing the extent to which the model lacks or possesses inherent planning capabilities. See Appendix A.1.3 for further details.

Results. As shown in Fig. 5, we observe that explicit task decomposition serves as a powerful intervention to improve LRM reasoning. Crucially, the greater efficacy of the multi-turn condition (gain > 10%) underscores the importance of isolating planning as a distinct cognitive step, suggesting that LRMs' intrinsic ability to plan is underdeveloped and requires explicit elicitation. This provides a firm empirical basis for our agentic framework, which is predicated on the principle that structured, upfront planning is a necessary precursor to reliable execution.

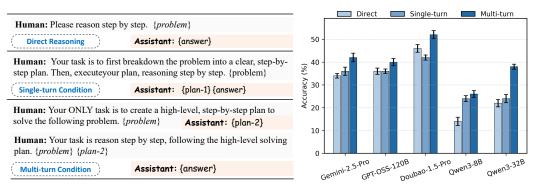


Figure 5: Validation of LRMs' intrinsic planning ability via task decomposition. We observe that explicit task decomposition enhances LRM reasoning, but the multi-turn setting is more effective. This pronounced gain validates the separation of planning and execution for reliable problem-solving.

4.4 METACOGNITIVE CONTROL: METACOGNITIVE FLEXIBILITY

We next measure metacognitive flexibility—the model's ability to adaptively shift reasoning strategies when recognizing that the current strategy isn't effective.

Experiment Details. We consider problems from DEEPMATH103K datasets, each augmented with three types of reasoning traps: *value corruption*, *unit corruption*, and *operation corruption*. In each case, an intermediate step is corrupted, and the CoT is truncated at that point. To evaluate whether models detect and adjust to the trap, we first ask models to continue reasoning. and compare the final answer accuracy against a baseline without corruption. We then test 5 cutting-edge models to judge whether the response correctly identifies and compensates for the corrupted step. The *flexibility*

rate is defined as the frequency with which a model successfully corrects the flawed reasoning and arrives at the correct solution. See Appendix A.1.4 for further details.

Results. We find that models often fail to recover from corrupted reasoning, continuing with invalid assumptions (i.e., Fig. 6). However, flexibility increases when the corruption is more obvious (e.g., extreme numerical distortions). These results reveal that metacognitive flexibility remains fragile and heavily reliant on superficial cues rather than deep structural awareness.

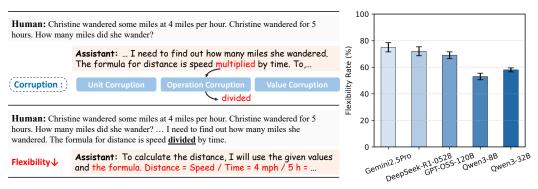


Figure 6: Validation of metacognitive flexibility. We find that models often fail to recover from corrupted reasoning, continuing with invalid assumptions, revealing fragile flexibility.

5 Towards Designing Metacognitive Reasoning Models

In §4, we demonstrate that current LRMs exhibit incomplete and fragile metacognitive abilities. Thus, we propose two paradigms for metacognitive enhancement: (1) **Emergent Metacognition**, an explicit prompting-based system for modular control, and (2) **Internalized Metacognition**, an intrinsically trained model for parameter-level metacognition.

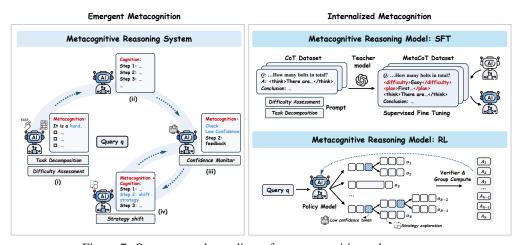


Figure 7: Our proposed paradigms for metacognitive enhancement.

5.1 THE PROMPT-DRIVEN METACOGNITIVE REASONING SYSTEM

First, we propose **Emergent Metacognition**, an explicit prompting framework designed to simulate a full metacognitive reasoning loop through modular API calls.

Experiment Details. The workflow, illustrated in Fig. 7, proceeds as follows: the model (i) self-assesses task difficulty and proposes a decomposition plan, (ii) executes the initial reasoning steps, (iii) dynamically identifies intermediate solutions with low confidence, and (iv) receives feedback and adaptively adjusts its strategy. The loop terminates when sufficient consistency is achieved (e.g., 5 consecutive verification passes) or persistent failure occurs (e.g., a 10-step failure streak). Each component is executed by the same underlying LRM architecture, instantiated independently

Table 2: Results of our Prompt-Driven Metacognitive Reasoning System. We evaluate the accuracy gain, and the flexibility rate increment on the corrupted DeepMath, as described in §4.4.

Models	Methods AIME2024 AIMI	AIME2025	GPQA	DeepMath (w/ corruption)	
		pass@1	pass@1	pass@1	flexibility rate
Gemini-2.5-Pro	Vanilla	90.8	83.0	83.0	75.2
	Ours	100.0	93.3	96.5	95.7
DeepSeek-R1-0528	Vanilla	91.4	87.5	81.0	71.9
	Ours	96.7	90.0	93.4	93.1

and prompted with a specific role aligned to a metacognitive function. This framework exposes latent metacognitive abilities such as self-awareness, task decomposition, confidence monitoring, and strategic flexibility by scaffolding higher-order control without requiring additional training. See Appendix A.2 for further details.

Results. We test this system on two strong models. Firstly, our system demonstrates significant improvements on 2 mathematical tasks and 1 QA benchmark. Both models show substantial performance gains (some even achieve 100%). Secondly, to further validate the efficacy of our metacognitive approach, we observed an increase of over 20% in flexibility rate on the corrupted DeepMath dataset. This enhancement effectively mitigates the flexibility deficit discussed in §4.4, underscoring the framework's ability to foster more adaptive and robust reasoning.

5.2 THE INTRINSIC METACOGNITIVE REASONING MODEL

While the prompt-driven system simulates metacognitive behavior through role-specific prompting, they do not endow the model with parameter-level metacognitive knowledge. To bridge this gap, we propose **Internalized Metacognition**, a intrinsic metacognitive reasoning model (MRM), which instills metacognitive functions through a two-stage approach: (1) supervised fine-tuning (SFT) as a cold start, followed by (2) reinforcement learning (RL).

Cold-start SFT. We first construct training data by augmenting samples from GSM8K and MATH with structured metacognitive traces: <difficulty> self-assessment, <plan> high-level decomposition, and <think> reasoning steps. These components are concatenated to form the full reasoning trajectories. These trajectory are then used to fine-tune models, enabling behaviors like self-evaluation and planning to emerge during inference.

RL. We build on the GRPO algorithm (Shao et al., 2024), which eliminates the value function and estimates the advantage in a group-relative manner. Formally, for each question q, GRPO samples a group of outputs $\{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot \mid q)$ and computes the token ratio $r_{i,t}(\theta) = \frac{\pi_{\theta}(o_{i,t}|q,o_{i,< t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|q,o_{i,< t})}$. It updates the policy by maximizing the objective:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}\left[\frac{1}{G}\sum_{i=1}^{G} \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left(\min\left(r_{i,t}(\theta) \, \hat{A}_{i,t}, \, \operatorname{clip}(r_{i,t}(\theta), 1-\varepsilon, 1+\varepsilon) \, \hat{A}_{i,t}\right) - \beta D_{\text{KL}}(\pi_{\theta} \| \pi_{\text{ref}})\right)\right], \quad (1)$$

where ε and β are hyper-parameters, and $\hat{A}_{i,t}$ is the group-normalized advantage.

Confidence Monitoring. To implement metacognitive monitoring during RL, we identify positions that exhibit low confidence along each rollout. We define token confidence at each position t as:

$$C_t = -\frac{1}{K} \sum_{v \in \text{Top-}K} \log \pi_{\theta}(v \,|\, q, o_{< t}). \tag{2}$$

A low-confidence position is detected at timestep t if $C_t \leq C$, where C is a predefined confidence threshold. This event indicates high uncertainty along the reasoning path.

Strategy Control. When a low-confidence state $s_t = (q, o_{\leq t})$ is detected, we fork the reasoning process. From this anchor state, we launch M new rollouts $\{o'^{(m)}\}_{m=1}^{M} \sim \pi_{\theta_{\text{old}}}(\cdot \mid s_t)$. Finally, all fully-formed trajectories—both the original G rollouts and all newly forked continuations—are collected into a single, unified batch for advantage calculation. Let this final batch of N trajectories be denoted by $\mathcal{B} = \{o_u\}_{u=1}^{N}$. The group-relative advantage is then computed across this entire dynamic set:

$$\hat{A}_u = \frac{r(o_u) - \text{mean}(\{r(o_v)\}_{v=1}^N)}{\text{std}(\{r(o_v)\}_{v=1}^N)},$$
(3)

Table 3: Performance of our intrinsic MRM on math reasoning and metacognitive tasks.

Methods	GSM8K	MATH500	AIME2024	DeepM	Iath
Nemous	Acc	Acc	Acc	difficulty assessment	flexibility rate
Qwen2.5-Math-7B	70.3	64.0	11.2	29.9	32.7
\hookrightarrow Ours: SFT (w/ difficulty)	79.1	75.4	13.3	60.8	36.0
\hookrightarrow Ours: SFT (<i>w</i> / <i>difficulty</i> + <i>plan</i>)	82.2	77.0	13.3	58.6	39.9
⇔GRPO	75.9	71.6	16.7	30.1	44.4
⊖Ours: RL	82.5	75.3	26.7	29.7	47.9
⇔Ours: SFT+RL	85.5	80.2	33.3	55.9	51.2

where $r(o_u) \in \{0,1\}$ is the outcome reward of trajectory o_u . The final objective becomes:

$$\mathcal{J}_{\text{ours}}(\theta) = \mathbb{E}_{o_u \in \mathcal{B}} \left[\frac{1}{|o_u|} \sum_{t=1}^{|o_u|} \min \left(r_{u,t}(\theta) \, \hat{A}_{u,t}, \, \operatorname{clip}(r_{u,t}(\theta), 1 - \varepsilon, 1 + \varepsilon) \, \hat{A}_{u,t} \right) - \beta D_{\text{KL}}(\pi_\theta \| \pi_{\text{ref}}) \right]. \tag{4}$$

Results. We validate our intrinsic training paradigm on <code>Qwen2.5-Math-7B</code> using curated subsets of GSM8K and MATH for training (details in Appendix A.3). As shown in Tab. 3, our methods demonstrate clear efficacy. Cold-start SFT on metacognitive traces significantly boosts performance. Notably, including difficulty assessments increases accuracy on that task by more than double, confirming that the model can internalize this monitoring skill. Our proposed method also substantially improves upon its GRPO baseline, especially on AIME accuracy and metacognitive flexibility. Crucially, combining SFT and RL yields the best overall performance, establishing new state-of-the-art results across the board and confirming that metacognitive abilities can be effectively internalized.

6 RELATED WORK

Understanding and Demonstrating metacognitive in LRMs. Recent efforts to enhance capabilities of LRMs have increasingly drawn inspiration from metacognition (Didolkar et al., 2024; Bilal et al., 2025; Wang et al., 2025a), which is the model's ability to monitor, evaluate, and control its own thought processes (Flavell, 1979; 1976). Early attempts have explored and demonstrated that metacognitive behaviors can be explicitly elicited through direct prompting to guide the models in self-reflection and self-evaluation (Wang & Zhao, 2023; Madaan et al., 2023; Liu et al., 2024). Despite these advances, existing works are highly rely on well-designed prompts, lacking the adaptability across diverse scenarios. In addition, they primarily focus on monitoring and evaluating the metacognitive abilities in LRMs, without involving any modifications to the model itself.

Instilling Metacognition in LRMs. There have been recent attempts to explore metacognition integration with LRMs, including training an external module to empower meta-thinking (i.e., Meta-Reasoner (Sui et al., 2025), MetaScale (Liu et al., 2025)) and exploring multi-agent systems to expand the intelligence boundary(i.e., ReMa (Wan et al., 2025), MPDF (Yang & Thomason, 2025)). While effective, these paradigms rely on external modules (e.g., inter-agent communication or meta thinker), rather than fostering an intrinsic faculty. In contrast, our work pursues a more holistic approach by introducing a functional framework. Through targeted training, both metacognitive knowledge and regulation are directly internalized into the model's parameters, enabling the development of a system that possesses metacognition as an autonomous, intrinsic capability rather than merely simulating it.

7 Conclusion

This work introduces a functional framework for metacognition in Large Reasoning Models (LRMs), distinguishing between internal metacognitive information and observable abilities. The authors empirically demonstrate that while LRMs possess predictive internal signals for reasoning outcomes, these signals do not consistently translate into effective monitoring or control behaviors. To address this gap, they propose two enhancement paradigms: a prompt-driven system that assigns modular metacognitive roles and an intrinsic training model that embeds these abilities directly into the LRM's parameters. Their findings suggest that integrating metacognitive reasoning improves task performance and offers a promising direction for future LRM development.

REFERENCES

- Rakefet Ackerman and Valerie A Thompson. Meta-reasoning: Shedding metacognitive light on reasoning research. In *International handbook of thinking and reasoning*, pp. 1–15. Routledge, 2017.
- Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K Arora, Yu Bai, Bowen Baker, Haiming Bao, et al. gpt-oss-120b & gpt-oss-20b model card. *arXiv* preprint arXiv:2508.10925, 2025.
- Roberto Araya. Do chains-of-thoughts of large language models suffer from hallucinations, cognitive biases, or phobias in bayesian reasoning? *arXiv preprint arXiv:2503.15268*, 2025.
- Ahsan Bilal, Muhammad Ahmed Mohsin, Muhammad Umer, Muhammad Awais Khan Bangash, and Muhammad Ali Jamshed. Meta-thinking in Ilms via multi-agent reinforcement learning: A survey. *arXiv preprint arXiv:2504.14520*, 2025.
- Melanie Cary and Lynne M Reder. Metacognition in strategy selection: Giving consciousness too much credit. In *Metacognition: Process, function and use*, pp. 63–77. Springer, 2002.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Peter Dayan. Metacognitive information theory. *Open Mind*, 7:392–411, 2023.
- Aniket Didolkar, Anirudh Goyal, Nan Rosemary Ke, Siyuan Guo, Michal Valko, Timothy Lillicrap, Danilo Jimenez Rezende, Yoshua Bengio, Michael C Mozer, and Sanjeev Arora. Metacognitive capabilities of llms: An exploration in mathematical problem solving. *Advances in Neural Information Processing Systems*, 37:19783–19812, 2024.
- Anastasia Efklides and Plousia Misailidi. *Trends and prospects in metacognition research*. Springer, 2010.
- Klaus Fiedler, Rakefet Ackerman, and Chiara Scarampi. Metacognition: Monitoring and controlling one's own knowledge, reasoning and decisions. *The psychology of human thought: An introduction*, pp. 89–111, 2019.
- John H Flavell. Metacognitive aspects of problem solving, 1976.
- John H Flavell. Metacognition and cognitive monitoring: A new area of cognitive—developmental inquiry. In *American psychologist*, volume 34, pp. 906. American Psychological Association, 1979.
- Yichao Fu, Xuewei Wang, Yuandong Tian, and Jiawei Zhao. Deep think with confidence. *arXiv* preprint arXiv:2508.15260, 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*, 2025a.
- Dong Guo, Faming Wu, Feida Zhu, Fuxing Leng, Guang Shi, Haobin Chen, Haoqi Fan, Jian Wang, Jianyu Jiang, Jiawei Wang, et al. Seed1. 5-vl technical report. *arXiv preprint arXiv:2505.07062*, 2025b.
- 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.
 - Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv* preprint arXiv:2103.03874, 2021.

- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec
 Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. arXiv
 preprint arXiv:2412.16720, 2024.
 - Doohyuk Jang, Yoonjeon Kim, Chanjae Park, Hyun Ryu, and Eunho Yang. Reasoning model is stubborn: Diagnosing instruction overriding in reasoning models. *arXiv preprint arXiv:2505.17225*, 2025.
 - Dancheng Liu, Amir Nassereldine, Ziming Yang, Chenhui Xu, Yuting Hu, Jiajie Li, Utkarsh Kumar, Changjae Lee, Ruiyang Qin, Yiyu Shi, et al. Large language models have intrinsic self-correction ability. *arXiv preprint arXiv:2406.15673*, 2024.
 - Qin Liu, Wenxuan Zhou, Nan Xu, James Y Huang, Fei Wang, Sheng Zhang, Hoifung Poon, and Muhao Chen. Metascale: Test-time scaling with evolving meta-thoughts. *arXiv* preprint *arXiv*:2503.13447, 2025.
 - Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36:46534–46594, 2023.
 - Thomas O Nelson and John Dunlosky. When people's judgments of learning (jols) are extremely accurate at predicting subsequent recall: The "delayed-jol effect". *Psychological Science*, 2(4): 267–271, 1991.
 - Elisabeth Norman, Gerit Pfuhl, Rannveig Grøm Sæle, Frode Svartdal, Torstein Låg, and Tove Irene Dahl. Metacognition in psychology. *Review of General Psychology*, 23(4):403–424, 2019.
 - Lirong Qiu, Jie Su, Yinmei Ni, Yang Bai, Xuesong Zhang, Xiaoli Li, and Xiaohong Wan. The neural system of metacognition accompanying decision-making in the prefrontal cortex. *PLoS biology*, 16(4):e2004037, 2018.
 - 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.
 - 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.
 - Yuan Sui, Yufei He, Tri Cao, Simeng Han, Yulin Chen, and Bryan Hooi. Meta-reasoner: Dynamic guidance for optimized inference-time reasoning in large language models. arXiv preprint arXiv:2502.19918, 2025.
 - Lev Tankelevitch, Viktor Kewenig, Auste Simkute, Ava Elizabeth Scott, Advait Sarkar, Abigail Sellen, and Sean Rintel. The metacognitive demands and opportunities of generative ai. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, pp. 1–24, 2024.
 - Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. *arXiv* preprint arXiv:2507.20534, 2025.
 - Ziyu Wan, Yunxiang Li, Xiaoyu Wen, Yan Song, Hanjing Wang, Linyi Yang, Mark Schmidt, Jun Wang, Weinan Zhang, Shuyue Hu, et al. Rema: Learning to meta-think for llms with multi-agent reinforcement learning. *arXiv preprint arXiv:2503.09501*, 2025.
 - Guoqing Wang, Wen Wu, Guangze Ye, Zhenxiao Cheng, Xi Chen, and Hong Zheng. Decoupling metacognition from cognition: A framework for quantifying metacognitive ability in llms. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 25353–25361, 2025a.
 - Qian Wang, Zhanzhi Lou, Zhenheng Tang, Nuo Chen, Xuandong Zhao, Wenxuan Zhang, Dawn Song, and Bingsheng He. Assessing judging bias in large reasoning models: An empirical study. *arXiv* preprint arXiv:2504.09946, 2025b.

Under review as a conference paper at ICLR 2026 Yuqing Wang and Yun Zhao. Metacognitive prompting improves understanding in large language models. arXiv preprint arXiv:2308.05342, 2023. 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. 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. Wei Yang and Jesse Thomason. Learning to deliberate: Meta-policy collaboration for agentic llms with multi-agent reinforcement learning. arXiv preprint arXiv:2509.03817, 2025. Nick Yeung and Christopher Summerfield. Metacognition in human decision-making: confidence and error monitoring. Philosophical Transactions of the Royal Society B: Biological Sciences, 367(1594):1310–1321, 2012.

A APPENDIX

A.1 FURTHER DETAILS AND RESULTS FOR §4

A.1.1 EVALUATING SELF-AWARENESS

This appendix provides additional details on the dataset construction and full experimental results for the metacognitive self-awareness task presented in §4.1.

Dataset Collection Our benchmark for the self-awareness task is constructed from the DEEP-MATH103K (He et al., 2025) dataset. To ensure a balanced evaluation across a wide spectrum of problem complexity, we randomly sampled a subset of 999 problems. These were then partitioned into three non-overlapping difficulty categories of 333 problems each, based on their official numerical ratings provided in the original dataset. The specific thresholds used for partitioning are as follows:

- Easy: 333 problems with a rating < 3.5.
- **Medium**: 333 problems with a rating between 3.5 and 6.5 (inclusive).
- **Hard**: 333 problems with a rating > 6.5.

Detailed Model Performance We present the detailed performance metrics for each evaluated model. For each model, we report the overall accuracy, the confusion matrix, and the per-category Precision, Recall, and F1-scores.

Gemini 2.5 Pro (Google). Achieved an overall accuracy of 58.5%. The model shows a tendency to misclassify Easy problems as Medium, indicating a potential conservative bias.

Table 4: Confusion Matrix for Gemini 2.5 Pro.

Table 5:	Per-Category	Metrics	tor	Gemini
2.5 Pro.				

Actual \ Pred.	Easy	Medium	Hard
Easy	147	176	10
Medium	25	230	78
Hard	16	110	207

Category	Precision	Recall	F1-Score
Easy	0.782	0.441	0.564
Medium	0.446	0.691	0.542
Hard	0.702	0.622	0.659

Seed-1.5-VL Pro. Achieved an overall accuracy of 61.6%. Similar to Gemini 2.5 Pro, it struggles with distinguishing Easy from Medium problems.

Table 6: Confusion Matrix for Seed-1.5-VL Pro.

Table 7: Per-	Category Metric	cs for Seed-1.5-
VL Pro.		

Actual \ Pred.	Easy	Medium	Hard
Easy	147	176	10
Medium	25	230	78
Hard	16	110	207

Category	Precision	Recall	F1-Score
Easy Medium	0.782 0.446	0.441 0.691	0.564 0.542
Hard	0.702	0.622	0.659

GPT-o3 (OpenAI). Achieved the highest overall accuracy of 69.0%. Its performance is more balanced across categories compared to other models, although it still shows some confusion between Medium and Hard problems.

Qwen-32B. This model exhibited a strong degenerative bias, classifying nearly all problems as Medium (98.3% of predictions). This resulted in high recall for the Medium category but near-zero recall for Easy and Hard, leading to a very low overall accuracy.

Table 8: Confusion Matrix for GPT-o3.

asy Medi	um Hard
76 22	8 29
	57 59 76 228 28 10

Table 10: Confusion Matrix for Qwen-32B.

Actual \ Pred.	Easy	Medium	Hard
Easy Medium Hard	3 1 0	321 330 331	9 2 2
	0		2

Table 9: Per-Category Metrics for GPT-o3.

Category	Precision	Recall	F1-Score
Easy	0.712	0.772	0.741
Medium	0.588	0.685	0.632
Hard	0.816	0.613	0.700

Table 11: Per-Category Metrics for Qwen-32B.

0.009	0.018 0.502 0.012

System Prompts:

You are an expert AI assistant specializing in mathematical reasoning. You possess advanced metacognitive capabilities. Your current task is to act as a "Problem Assessor". Given a mathematical problem, your goal is to analyze its requirements and assess its difficulty for an AI like yourself. Do NOT solve the problem. You must only provide your assessment.

User Prompts:

Here is the problem: {problem_text}

Your task is to assess the difficulty of a mathematical problem based on the provided rubric and examples.

Difficulty Rubric

Easy: The problem follows a single, linear computational path using a standard formula or definition. The solution is straightforward and requires no creative insight.
Medium: The problem requires a sequential composition of

distinct conceptual modules or formulas. The solution involves a multi-step, but generally standard, reasoning process.

Hard: The problem requires a non-linear or exploratory reasoning path. The solution may demand non-obvious insights, creative problem transformations, or the synthesis of concepts from different mathematical branches.

Provide your response as a single JSON object wrapped in a markdown code block. The JSON object must contain the following keys:

"Difficulty_category": string, choose one from ["Easy", "Medium", "Hard"].

"Rationale": string, a brief explanation for your choice, explicitly referencing the rubric criteria.

Your entire output must be in the following format:

```
'`'json
{
"Difficulty_category": "...",
"Rationale": "..."
}
```

A.1.2 EVALUATING CONFIDENCE

This appendix provides formal definitions for the confidence metrics and further details on the experimental setup for the metacognitive confidence adjustment task presented in §4.2.

Confidence Metric Definitions

Token Confidence. Following standard practice (Fu et al., 2025), we define our base metric, token confidence, at each position t of a reasoning trace. It is calculated as the negative average log-probability of the top-k most likely tokens in the softmax distribution at that step:

$$C_t = -\frac{1}{k} \sum_{i=1}^{k} \log P(token_j \mid o_{< t}),$$
 (5)

where $P(token_j \mid o_{< t})$ is the probability of the j-th most likely token given the preceding sequence $o_{< t}$. Lower values of C_t correspond to higher model confidence (a more peaked distribution). For all our experiments, we set k=20.

Average Trace Confidence. To obtain a single confidence score for an entire reasoning trace of length N, we compute the average trace confidence by averaging the token confidences across all generated tokens:

$$C_{\text{avg}} = \frac{1}{N} \sum_{t=1}^{N} C_t.$$
 (6)

While useful as a global measure, C_{avg} can obscure critical, localized moments of uncertainty within a long reasoning process.

Evaluating Confidence Dynamics and Adjustment.

To analyze the model's confidence adjustment, as described in the main text, we introduce metrics designed to capture both the trajectory and the weakest points of a model's confidence.

Segmented Trace Confidence. To analyze the trajectory of confidence, we partition each reasoning trace into three equal, non-overlapping segments: **Start, Middle**, and **End**. We then compute the average trace confidence independently for each segment. These three scores, $(C_{\text{start}}, C_{\text{middle}}, C_{\text{end}})$, form the basis for our Confidence Trajectory Analysis. A trace is classified as *Increasing* if C_{end} is significantly lower (i.e., more confident) than C_{start} , and vice-versa for *Decreasing*. The *Consistently High/Low* patterns are determined by comparing all three segment scores against a predefined threshold.

A.1.3 EVALUATING TASK DECOMPOSITION

Task decomposition Prompt 1.

```
You are a helpful assistant. Solve the following mathematical problem. Please reason step by step and provide your final answer within \begin{tabular}{l} boxed{}. \end{tabular}
```

```
Problem:
---
{problem_text}
---
```

Task decomposition Prompt 2.

You are a helpful assistant. Your task is to first break down the problem into a clear, step-by-step plan. Then, execute your plan, reasoning step by step. Finally, provide your final answer within \boxed{}.

```
boxed{}.
Problem:
---
{problem_text}
---
```

Task decomposition Prompt 3. You are a meticulous problem-solving planner. Your ONLY task is to create a high-level, step-by-step plan to solve the following mathematical problem. The plan should consist of concrete, actionable steps. Do NOT actually solve the problem or perform any calculations. Problem: {problem_text} Plan: You are an expert problem solver. You will be given a problem and a pre-made plan. Your task is to follow this plan meticulously to solve the problem. Please reason step by step based on the plan. Finally, provide your final answer within \boxed{}. Problem: {problem_text} Plan: generated_plan_from_3a Solution:

864 A.1.4 EVALUATING METACOGNITIVE FLEXIBILITY 865 866 867 868 869 870 871 872 873 System Prompts: 874 You are a data generator for reasoning robustness evaluation. 875 Your task is to take an original math/logic problem with its 876 reasoning process (chain of thought) and answer, and then 877 intervene in exactly ONE reasoning step with a corruption. The 878 corruption must be *critical enough to change the final answer*. 879 880 User Prompts: 881 Input 882 Question: 883 {original_question} 884 Original Reasoning Process: 885 {reasoning_steps} 886 887 Answer: 888 889 {answer} 890 891 INSTRUCTION 892 1. Choose exactly ONE corruption type from the following three 893 (do not mix): 894 - Value Corruption: replace a key number with an incorrect but 895 plausible value (e.g., change 1 \rightarrow 1.111111 or g=9.8 \rightarrow g=1000). 896 - Unit Corruption: replace the unit of a key step with another 897 898 (e.g., meters \leftrightarrow centimeters, hours \leftrightarrow minutes). 899 - Operation Corruption: change the mathematicallogical operation 900 in a key step 901 (e.g., replace + with -, union with intersection, inequality with 902 equality). 903 904 2. Apply the corruption to ONE critical reasoning step. 905 - Make sure the corruption influences the correctness of the final 906 answer. - After the corruption, truncate the reasoning at that corrupted 907 step (do not continue to the correct answer). 908 909 3. Produce the final output strictly in the following JSON 910 format: 911 912 "question": "...", 913 "corrupted_reasoning": "...", // include the question + steps up 914 to the corrupted step

"corruption_type": "Value Corruption | Unit Corruption | Operation

915

916

917

Corruption"

}

```
918
      You are a reasoning evaluator.
919
      Your task is to judge whether a model successfully recognized and
      adjusted to a corrupted reasoning trap. You are given:
921
      - The original problem
922
      - The ground truth answer
      - The type of corruption applied
923
      - The last corrupted reasoning step (truncated point)
      - The model's continued reasoning and final answer
925
926
      You must decide: Did the model detect and flexibly adjust to the
927
      trap?
928
929
930
      System Prompts:
      You are a reasoning assistant. You are given a math/logic
931
      problem, together with a partially completed reasoning process.
932
      Please continue reasoning from what is provided.
933
934
      User Prompts:
935
      INPUT
936
      Question: {original_question}
937
938
      Ground Truth Answer: {ground_truth_answer}
939
      Ground Truth Corruption Type: {ground_truth_corruption_type}
940
      Last Corrupted Reasoning Step: {last_corrupted_step}
941
      Model Generated Reasoning:
942
      model_generated_cot
943
944
      Model Final Answer:
945
      model_final_answer
946
947
      TASK Based on the evidence:
948
      - If the model explicitly or implicitly identified the corrupted
949
      reasoning (e.g., points out error, discards it, corrects it)
      and produced the correct final answer, output "Yes".
951
      - Otherwise (if it followed the trap blindly, failed to adjust, or
952
      produced the wrong answer), output "No".
953
      Final Output (strict format):
954
      Recognition And Adjustment: Yes | No
955
```

Prompt of Step 1

956

```
972
      System Prompt:
973
      You are an elite mathematical strategist and analyst. Your
      primary function is to perform a deep Metacognitive Analysis of
975
      complex mathematical problems. You are to deconstruct the problem
      into its core components, identify underlying principles, and then
976
      formulate a high-level, executable strategic plan.
977
      Your task is to produce a Metacognitive Analysis of the following
978
      problem. You must NOT provide a final solution or perform
979
      detailed calculations.
980
      ### Core Principles
981
      * **Analytical Depth:** Your analysis must go beyond a
982
      surface-level reading. Identify the mathematical field, key
983
      concepts, constraints, and the explicit goal.
984
      * **Strategic Foresight:** Your plan should be a viable path to
985
      a solution. This includes anticipating potential difficulties,
986
      identifying necessary lemmas, and choosing the most promising
987
      approach.
      * **Clarity and Brevity: ** The analysis and plan must be clear,
988
      concise, and easily understood by another mathematical expert who
989
      will execute it.
990
991
      User Prompt:
992
      ### Your Task
993
      **Problem:**
994
      _____
995
      {problem_statement}
996
      ========
997
      **Metacognitive Analysis:**
      **1. Problem Deconstruction:**
998
      * **Mathematical Domain:** Identify the primary field(s) of
999
      mathematics involved (e.g., Number Theory, Combinatorics,
1000
      Euclidean Geometry).
1001
      * **Given Conditions & Constraints:** List all the premises,
1002
      conditions, and constraints provided in the problem statement in
1003
      a structured format.
      * **Objective: ** State the precise question to be answered or the
1005
      proposition to be proven.
      **2. Strategic Solution Plan (Method Sketch): **
1007
      Present a high-level, conceptual outline of your proposed solution
      path. This sketch should enable an expert to grasp the entire
      logical flow of the argument without needing the full details.
1009
      must include:
1010
      * **Overall Strategy Narrative:** A brief description of the core
1011
      idea behind your approach (e.g., "We will use proof by induction,"
1012
      "The strategy is to establish a coordinate system and use analytic
1013
      geometry, " "We will prove the contrapositive by assuming...").
1014
      * **Key Lemmas and Intermediate Results:** State the full and
1015
      precise mathematical formulations of any key lemmas or theorems
1016
      you plan to prove or apply. These are the major milestones of the
1017
      proof.
1018
      * **Logical Skeleton: ** If applicable, describe the key
      constructions, case splits, or transformations that form the
1019
1020
      backbone of your argument.
      * **Potential Challenges & Pitfalls: ** Briefly note any steps
1021
      that might be particularly tricky, prone to error, or require a
1022
      non-obvious insight.
1023
      ### Negative Constraints
```

```
1026
      * **DO NOT** write the full, step-by-step solution.
1027
      * **DO NOT** perform detailed algebraic manipulations or numerical
      calculations.
1029
      * Your output should be strictly limited to the analysis and
1030
      strategic plan as outlined above.
1031
      Prompt of Step 2
1032
1033
      System Prompt:
1034
      You are an exceptionally rigorous mathematical solver.
      sole purpose is to take a pre-defined strategic plan and execute
1035
      it with absolute precision and logical soundness. You must not
1036
      deviate from, question, or reinterpret the provided plan.
1037
      Your task is to produce a complete and formally justified solution
1038
      to the following mathematical problem, strictly following the
1039
      'Solution Plan'.
1040
      ### Core Principles
      * **Rigor is Paramount:** Your primary goal is to produce a
1042
      complete and rigorously justified solution. Every step in your
1043
      solution must be logically sound and clearly explained. A correct
1044
      final answer derived from flawed or incomplete reasoning is
      considered a failure.
      * **Unyielding Adherence to Plan: ** You MUST strictly follow
1046
      the logical flow, lemmas, and constructions laid out in the
      'Solution Plan'. Do not introduce new methods, skip steps, or
1048
      alter the proposed strategy in any way. Your role is execution,
1049
      not creation.
1050
      * **Honesty About Completeness:** If you cannot find a complete
1051
      solution following the plan, you must **not** guess or create a
1052
      solution that appears correct but contains hidden flaws. Instead,
1053
      you should present only the significant partial results that you
1054
      can rigorously prove by following the plan.
1055
1056
      User Prompt:
      ### Your Task
1057
      **Problem:**
1058
      _____
1059
      {problem_statement}
      ========
1061
      **Solution Plan:**
1062
      ========
1063
      {step1_output}
1064
1065
      **Detailed Solution: **
1066
      Present the full, step-by-step mathematical proof, meticulously
      following the guidance of the 'Solution Plan'. Each step must be
1067
      logically justified and clearly explained. The level of detail
1068
      should be sufficient for an expert to verify the correctness of
1069
      your reasoning without needing to fill in any gaps. This section
1070
      must contain ONLY the complete, rigorous proof, free of any
1071
      internal commentary, alternative approaches, or failed attempts.
1072
      * **Use TeX for All Mathematics:** All mathematical variables,
1073
      expressions, and relations must be enclosed in TeX delimiters
1074
      (e.g., 'Let $n$ be an integer.').
1075
```

Your step-by-step reasoning, strictly following the plan, begins

After completing the detailed solution, state the final answer

1076

1077

1078

1079

here...

Final Answer

within \boxed{}.

1080 **Prompt of Step 3** 1081 System Prompt: 1082 You are an expert mathematician and a meticulous grader for 1083 AIME-level computational problems. Your primary task is to 1084 rigorously verify the provided solution's **computational 1085 reasoning and numeric correctness**. A solution is to be judged 1086 correct **only if every step that affects the numeric outcome 1087 is correct and sufficiently justified.** A solution that reaches 1088 a correct final integer answer via arithmetic slips, incorrect 1089 algebraic manipulations, unverified casework, counting mistakes, 1090 or hidden assumptions must be flagged as incorrect or incomplete. 1091 ### Instructions ### 1092 **1. Core Instructions** * Your sole task is to identify and report all issues in the 1093 provided solution. You must act strictly as a **verifier**, NOT 1094 a solver. 1095 * You must **NOT attempt to correct, fix, or complete** any errors 1096 or missing arguments. 1097 * Perform a **step-by-step** check of the entire solution and 1098 produce a **Detailed Verification Log**. For each step: 1099 * If the step is correct, state briefly that it is correct. 1100 * If the step contains an issue, explain the error and classify it 1101 (see section 2). 1102 **2. How to Handle Issues in the Solution** 1103 All issues must be classified into one of the following categories: 1104 * **a. Critical Error:** 1105 * Definition: Any error that changes or potentially invalidates 1106 the numeric result. Examples include arithmetic mistakes, 1107 wrong algebraic transformations, misapplied formulas, incorrect 1108 combinatorial counts, invalid casework, or unjustified 1109 approximations that affect the integer outcome. 1110 * **Procedure:** 1111 * Point out the exact error and explain why it invalidates the reasoning. 1113 * Do **not** check further steps that rely on this error. 1114 * You may still check other independent parts of the solution. * **b. Justification Gap:*** Definition: Steps where the stated 1115 conclusion might be correct, but the reasoning is incomplete or 1116 not justified at AIME level. 1117 * **Procedure:** 1118 * Point out the missing justification. 1119 * Explicitly state that you will assume the step's conclusion 1120 holds for the sake of checking subsequent steps. 1121 **3. Output Format** 1122 Your response MUST be structured into two main sections: a 1123 **Summary** followed by the **Detailed Verification Log**. 1124 * **a. Summary*** **Final Verdict:** One clear sentence declaring overall validity (e.g., "The solution is correct," "The solution 1125 contains a Critical Error and is therefore invalid," or "The 1126 solution contains several Justification Gaps."). 1127 * **List of Findings:** A bulleted list of every issue found. 1128 each finding include: 1129 * **Location:** A direct quote of the key phrase or equation. 1130

21

* **Issue:** Short description and classification (**Critical

* **b. Detailed Verification Log*** Provide a step-by-step

Error** or **Justification Gap**).

verification.

1132

```
1134
      * Quote the relevant part of the solution before your check.
1135
      * State clearly: **Correct**, **Critical Error**, or
      **Justification Gap**.
1137
      * Do **not** supply corrections or alternative methods | only
1138
      report the issues.
      **Important:**
1139
      - Do not propose fixes or alternative solutions.
1140
      - Do not attempt to supply missing reasoning.
1141
      - Only check and report correctness of what is written.
1142
1143
      User Prompt:
1144
      ### Your Task
1145
      **Original Problem:**
1146
      ========
1147
      {problem_statement}
1148
      ========
1149
      **Current Solution:**
      _____
1150
      {last_solution}
1151
      _____
1152
1153
      ### Monitoring Task Reminder ###
1154
      Your task is to act as an math grader. Now, generate the
1155
      **summary** and the **step-by-step verification log** for the
1156
      solution above. In your log, justify each correct step and
1157
      explain in detail any errors or justification gaps you find, as
1158
      specified in the instructions above."
1159
```

Prompt of Step 4

```
1188
      System Prompt:
1189
      You are an expert mathematician and a careful corrector for
      AIME-level computational problems.
1191
      You will be given three inputs:
      1) The Original Problem,
1192
      2) The current Solution,
1193
      3) A Verification Log (from a previous check), which labels each
1194
      step as Correct / Justification Gap / Critical Error, and provides
1195
      short notes.
1196
      ### Your Task ###
1197
      Using the Verification Log, **step by step correct the Original
1198
      Solution**.
1199
      - If a step is labeled **Correct**, keep it unchanged (you may
1200
      lightly reformat for clarity).
1201
      - If a step is labeled **Justification Gap**, supply the missing
1202
      justification or intermediate calculations, enough for AIME-level
1203
      rigor.
      - If a step is labeled **Critical Error**, replace it with
1204
      a correct mathematical step (with explicit computations or
1205
      reasoning) and update all dependent later steps accordingly.
1206
      - Do **not** introduce new solution paths, alternative methods, or
1207
      multiple approaches. Only repair the given solution chain.
1208
      ### Output Format ###
1209
      1. **Correction Summary**
1210
      - A single sentence declaring whether the solution has been fully
1211
      corrected and what the final answer is.
1212
      - Example: \The solution has been fully corrected. Final Answer
1213
      = 70."
      - Or, if not possible: \The solution cannot be fully corrected
1214
      due to missing information in step X."
1215
      2. **Correction Log**
1216
      For each relevant step (especially those flagged in the
1217
      Verification Log), provide an entry with:
1218
      - **Quoted Step: ** The original line/equation (quoted or in a code
1219
      block).
1220
      - **Verification Label:** Correct / Justification Gap / Critical
1221
      Error.
1222
      - **Correction / Action:**
1223
      * If Correct → \Unchanged | correct."
1224
      * If Justification Gap → Provide the missing computation/derivation
      briefly, ending with \Filled gap."
1225
      * If Critical Error → Provide the corrected computation/derivation,
1226
      briefly note why the original was wrong, and end with \Corrected."
1227
      - If a step's correction affects later steps, explicitly note
1228
      \Affects subsequent steps: Yes/No."
1229
      3. **Full Corrected Solution**
1230
      - Present the entire solution in a clean, continuous write-up,
1231
      combining unchanged and corrected steps.
1232
      - Show all necessary algebra, arithmetic, or combinatorial
1233
      reasoning clearly.
1234
      - After completing the detailed solution, state the final answer
1235
      within \boxed{}.
1236
```

```
1242
       User Prompt:
1243
      ### Your Task
      ***Original Problem:**
      _____
1245
      {problem_statement}
1246
      _____
1247
      **Current Solution:**
1248
      ========
1249
      {last_solution}
1250
      ========
1251
      **Verification Log: **
1252
      ========
1253
      {monitor_output}
1254
      _____
```

Stability and efficiency refinements. (i) Dynamic sampling. We reject prompts whose group rewards are all 0 or all 1 (no learning signal), and resample until a batch with informative gradients is formed—stabilizing updates as training progresses. (ii) Token-level aggregation. We aggregate losses at the token level rather than per-sample, ensuring longer sequences contribute proportionally and preventing degenerate length/entropy dynamics in long-CoT training. (iii) Optional KL. In pure reasoning RL we set β =0, removing the KL constraint as the optimal policy may significantly diverge from the reference model.

A.2 STATEMENT ON THE USE OF LARGE LANGUAGE MODELS (LLMS)

LLMs as the Subject of Research. One of the core components of our research involves the investigation and evaluation of the reasoning capabilities of current open-source and closed-source Large Language Models (LLMs). As such, a number of LLMs are explicitly named and analyzed within this paper (as detailed in §4). In this capacity, they serve as the objects of our study.

LLMs as an Assitive Tool. In the preparation of this manuscript, the use of LLMs is limited to polishing the text for grammatical correctness, spelling, and clarity of expression. The LLMs were not used to generate any core research ideas, experimental designs, data analysis, or substantive portions of the manuscript.

We assume full responsibility for all content presented in this paper, including any text that has been revised with the assistance of an LLM. We have meticulously reviewed and edited all content to ensure its scientific accuracy and originality, preventing any form of plagiarism or academic misconduct.

A.3 IMPLEMENTATION DETAILS

For the online model APIs, we utilized their respective official endpoints with default temperature settings. The maximum generation length was also kept at its default value, and no other parameters were modified.

For the training of open-source models, all experiments were conducted on a server equipped with 8×NVIDIA H800 GPU (80GB). The maximum generation length was set to 8192 for the GSM8K and MATH datasets, and 16384 for all other datasets. The temperature was set to 1.0, and the random seed was fixed to 42 for reproducibility.

For SFT training, we adpot full-parameter fine-tuning with learning rate $r=1\times 10^{-5}$ and batchsize=32. To construct training data, we first define the difficulty level for each sample from the train set of GSM8K and MATH: Easy (GSM8K and MATH lv.1), Medium (MATH lv.1-4), and Hard (MATH lv.3-5). Second, to obtain high-level task decomposition, we utilize Gemini-2.5-Pro with the prompt in Appendix A.1.3. The final pattern of training samples in cold-start is in the form of: "<difficulty> level </difficulty> <plan> decomposition </plan> COT

For RL training, we write our code based on the open-source Verl framework. Training settings are listed in Tab. 12.

Table 12: Training settings for RL.

Parameter	Value
n_gpu	8
rollout.n	16
total_steps	1000
batch_size	8
critic_warmup	0
max_prompt_length	512
max_response_length	16384
filter_overlong_prompts	True
learning_rate	1e-6
use_kl_loss	True
kl_loss_coef	0.001
kl_loss_type	low_var_kl