

RiT: Rubrics-in-Thinking Reinforcement Learning for Improved Reasoning in Large Language Models

Anonymous ACL submission

Abstract

Large Reasoning Models (LRMs) benefit from generating intermediate reasoning steps, enabling more reliable and interpretable decision-making. While outcome-based supervision has proven effective for LRMs across diverse tasks, it focuses solely on final answers and cannot guarantee high-quality intermediate reasoning. In contrast, existing process supervision is largely limited to verifiable domains such as mathematics or code, where intermediate steps can be explicitly checked, restricting its applicability to open-ended reasoning tasks. To address these limitations, we propose Rubrics-in-Thinking Reinforcement Learning (RiT), the first framework to introduce *thinking-rubric supervision* into intermediate reasoning. RiT automatically generates fine-grained rubrics and integrates them into a reward function via gated fusion with outcome-based rewards, guiding models to reason in a coherent and task-aligned manner, improving both intermediate steps and the final response. Experiments on reasoning-intensive and open-ended benchmarks demonstrate that RiT consistently outperforms outcome-only RL baselines.

1 Introduction

Large reasoning models (LRMs) (DeepSeek-AI, 2025; Jaech et al., 2024; Wei et al., 2022) have recently attracted attention for generating intermediate reasoning steps alongside final answers, enabling more reliable and interpretable decision-making. This capability benefits diverse natural language tasks, including mathematical problem solving (Mai et al., 2025), multi-step question answering (Lee and Hockenmaier, 2025), logical inference (Chen et al., 2025), and code generation (El-Kishky et al., 2025). However, most training paradigms for LRMs rely heavily on outcome-based rewards, optimizing primarily for final answer correctness. This can encourage shortcut reasoning, break multi-step logical consistency,

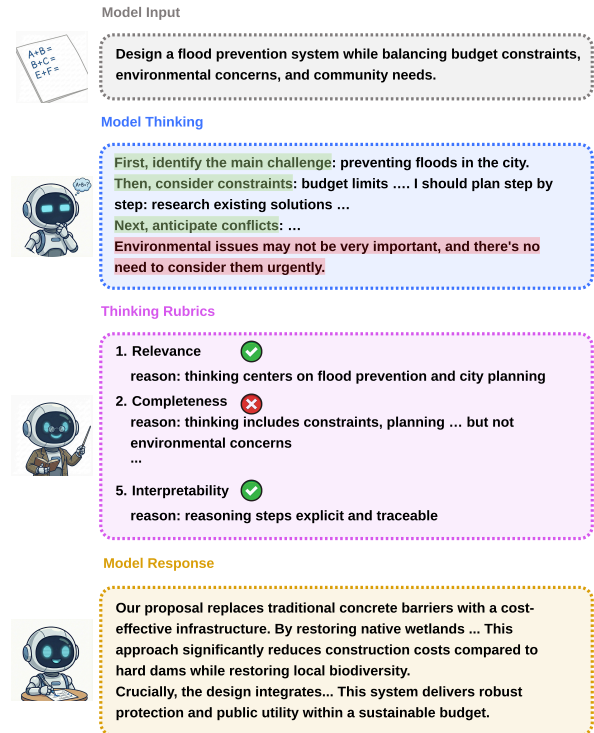


Figure 1: Thinking rubrics guide intermediate reasoning to maintain quality and task alignment.

and overlook qualitative aspects such as coherence and decisiveness. Process reward models (PRMs) (Khalifa et al., 2025) offer one solution but are often domain-specific (Lu et al., 2025; Fan et al., 2025; Yu et al., 2025b), relying on hand-crafted heuristics or verifiable intermediate steps, and thus struggle to generalize to open-ended tasks involving subjective judgment or semantic interpretation. These limitations highlight a key challenge: **how can we guide a model's reasoning in a flexible yet principled way?**

One promising direction comes from rubric-based reinforcement learning (RL), which has recently shown success in guiding model behavior through structured criteria for evaluating final responses (Gunjal et al., 2025; Liu et al., 2025; Huang

et al., 2025). However, the application of rubrics to intermediate reasoning remains largely unexplored. As Figure 1 illustrates, extending rubric evaluation to the thinking process enables interpretable, multi-dimensional feedback on reasoning trajectories, improving both quality and reliability. Motivated by this, we propose Rubrics-in-Thinking RL (RiT), a simple and general framework that automatically generates fine-grained thinking rubrics and integrates them into a reward function for RL. This function evaluates reasoning trajectories according to rubrics and is fused with outcome-based rewards to balance the quality of reasoning and final response, encouraging more coherent and interpretable reasoning across tasks.

Our main contributions are:

- We systematically study thinking-rubric supervision in LRMs, analyzing its effects.
- We propose RiT, which automatically generates fine-grained rubrics and leverages them to supervise intermediate reasoning via RL.
- Experiments demonstrate the effectiveness and generalization of thinking-rubric supervision across diverse tasks and models.

2 Related Work

Large reasoning models (LRMs) (DeepSeek-AI, 2025; Marjanovic et al., 2025; Jaech et al., 2024) are a class of large language models (LLMs) designed to produce intermediate reasoning steps alongside final responses. By explicitly exposing reasoning chains, LRMs enable verification, planning, and reasoning-aware optimization, distinguishing them from conventional LLMs that output only end results (Chen et al., 2025; Xu et al., 2025). Existing LRMs have been primarily developed for narrow domains such as mathematical problem solving, code generation, and formal logic, often leveraging reward models or reinforcement learning to supervise the coherence and correctness of reasoning steps (Mai et al., 2025; El-Kishky et al., 2025; Fan et al., 2025; Lee and Hockenmaier, 2025; Khalifa et al., 2025; Yu et al., 2025b). However, these approaches rarely generalize to heterogeneous tasks, where inputs, outputs, and evaluation metrics vary significantly across domains, such as open-ended generation and general reasoning benchmarks. To address this limitation, we propose a thinking-rubric-based supervision framework that guides LRM training across diverse tasks,

promoting reasoning-aware behavior beyond specialized domains.

Rubric-based RL has emerged as a paradigm for aligning LLMs by replacing scalar rewards with structured, multi-dimensional evaluation criteria. While traditional RL with Verifiable Rewards (RLVR) (Shao et al., 2024; Yu et al., 2025a; Zheng et al., 2025) excels in objective domains like math and code, recent works (Gunjal et al., 2025; Huang et al., 2025) show that rubrics can extend RL to open-ended reasoning tasks without ground truth. To reduce manual design, methods such as OpenRubrics (Liu et al., 2025) enable automated, scalable rubric generation, capturing nuanced human preferences. Rubrics further aid exploration and instruction following (Zhou et al., 2025), and have been specialized for high-stakes or subjective domains (Wang et al., 2025b), collectively advancing interpretable and robust reward modeling in RL for LLM. Notably, however, no prior work has systematically investigated the specific impact of thinking rubrics on model reasoning performance.

Open-ended generation tasks (Que et al., 2024; Bai et al., 2025; Lin et al., 2025) such as long-form and professional writing have become a primary focus of recent research as model capabilities advance beyond traditional closed-form settings. These tasks lack a single canonical target and instead require satisfying multiple qualitative criteria. Traditional metrics are often inadequate (Liu et al., 2023b), motivating the adoption of LLM-as-a-judge frameworks to assess generation quality (Gu et al., 2024; Bai et al., 2024). Methodologically, prior work on open-ended generation primarily adopts either data-centric enhancement or outcome-oriented reinforcement learning (Pham et al., 2024; Jia et al., 2025; Wang et al., 2025a).

3 Preliminaries

Formalization of LRMs. Given an input query q , an LRM generates an output sequence that consists of an explicit reasoning component and a final answer. Formally, the output is defined as:

$$y = \langle r, o \rangle, \quad (1)$$

where r denotes a sequence of intermediate reasoning steps and o denotes the final response. r is explicitly delimited by a special structural token `</think>`, which separates the internal reasoning process from the final answer. In this work, we refer to r as the *thinking process* of the LRM.

Reinforcement Learning for LRMs. We define $\pi_\theta(y | q)$ as a parameterized policy mapping an input query q to an output sequence y . Reinforcement learning is employed to optimize the policy with respect to a scalar reward function $R(q, y)$, which evaluates the quality of the generated output, including both the final response and the associated reasoning process. Specifically, the learning objective is defined as:

$$\max_{\theta} \mathbb{E}_{y \sim \pi_\theta(\cdot | q)} [R(q, y)]. \quad (2)$$

4 Methodology

4.1 Overview

In this section, we present our proposed RiT framework. Our key insight is that effective reasoning in LRMs necessitates the explicit regulation of intermediate thinking processes, rather than merely optimizing the final response quality. To this end, we propose a unified training framework that integrates thinking-aware reward modeling with group-relative optimization, enabling stable and interpretable learning under heterogeneous reward signals. The technical implementation of RiT consists of three key components, as illustrated in Figure 2. First, we construct a set of **Thinking Rubrics** that specify task-dependent criteria for evaluating thinking processes, explicitly defining what constitutes desirable reasoning behaviors under different queries. Second, based on these rubrics, we design a **Gated Fusion Reward** that reflects both thinking and response quality while mitigating spurious high rewards from superficial patterns. Finally, we optimize the model using **Group Relative Policy Optimization (GRPO)** (Shao et al., 2024), where the fused reward serves as the training signal and group-wise normalization enables stable learning across heterogeneous reward scales.

4.2 Thinking Rubrics Generation

The first step in RiT is to generate **thinking rubrics**, i.e., explicit criteria that guide and evaluate the model’s intermediate reasoning. Reasoning in LRMs is often under-specified or unstructured, which can result in incomplete, inconsistent, or logically flawed trajectories. By providing these rubrics, we establish clear standards for high-quality reasoning, enabling reliable evaluation of intermediate outputs and effectively guiding subsequent RL. To generate the rubrics, a frontier LLM is prompted with instructions comprising:

- **Task Definition.** The prompt specifies the task requirements and, when applicable, the evaluation dimensions of the final response (e.g., factual correctness, relevance, completeness). This ensures that the generated rubrics are aligned with the task objectives and provide explicit guidance for reasoning.
- **Positive Examples.** Model-generated reasoning trajectories that successfully satisfy the task requirements. These exemplify desirable reasoning behaviors and facilitate the LLM in inferring high-quality rubrics.
- **Negative Examples.** Reasoning trajectories that fail to meet the task requirements. These illustrate common errors, enabling the LLM to generate rubrics that effectively distinguish correct from incorrect reasoning.

After generating an initial rubric set, we further employ the LLM to iteratively refine the rubrics through a self-evaluation process. Specifically, the LLM is instructed to examine whether the current rubrics suffer from *missing criteria*, *redundant or overlapping conditions*, or *ambiguous*. Based on this analysis, the LLM revises the rubrics by adding missing aspects, removing or merging suboptimal criteria, and improving clarity and coverage.

The resulting thinking rubrics are denoted as:

$$\mathcal{C}_{\text{think}} = \{c_1, c_2, \dots, c_K\}, \quad (3)$$

where each c_k represents an explicit criterion that defines a desirable aspect of reasoning (e.g., logical coherence, relevance, completeness). These rubrics form the foundation for computing the thinking reward $R_{\text{thinking}}(r)$, ensuring that the model’s reasoning trajectories are structured, task-aligned, and interpretable. Importantly, this design is flexible and general, applicable to a wide range of tasks including open-ended generation and structured reasoning without task-specific modifications.

4.3 Reward Design

Thinking Rubric Reward. This reward quantifies the quality of the model’s intermediate reasoning, explicitly incentivizing trajectories that are well-formed, coherent, and aligned with task-specific criteria. High-quality reasoning not only improves the final response but also enhances interpretability and controllability of model behavior.

For a reasoning trajectory r corresponding to query q , its performance is evaluated against the

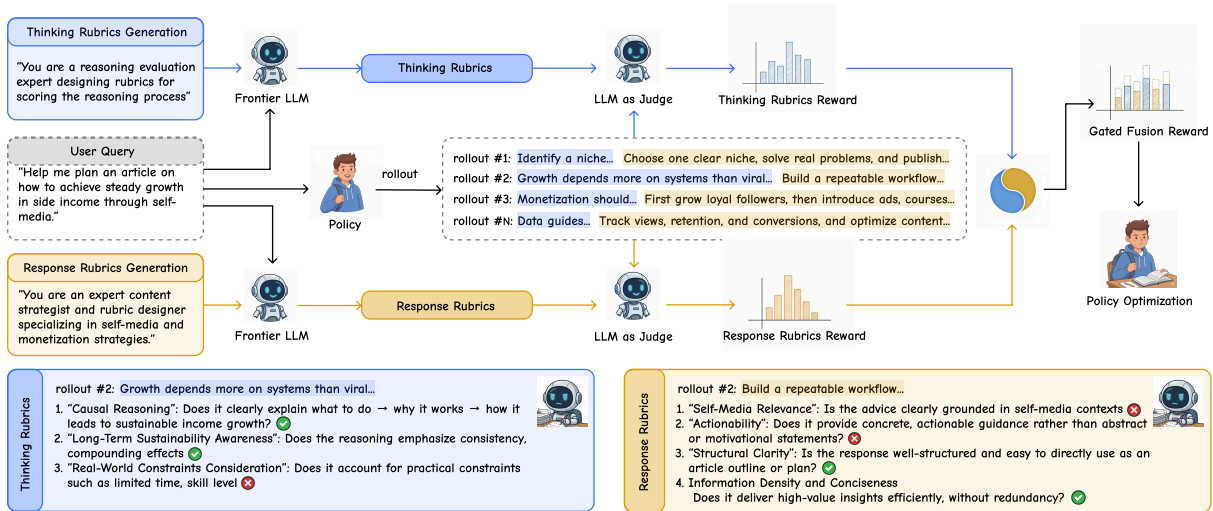


Figure 2: Overview of RiT, where a thinking-rubric reward is integrated with the response reward via gated fusion.

rubrics set $\mathcal{C}_{\text{think}}$. Formally, the thinking rubric reward is defined as:

$$R_{\text{thinking}}(r) = \frac{1}{K} \sum_{k=1}^K s_k^{\text{think}}, \quad (4)$$

where $s_k^{\text{think}} \in \{0, 1\}$ indicates whether r satisfies the k -th criterion under the LLM-based evaluator. By leveraging a frontier LLM to score the rubrics, this design ensures evaluations that are both robust and adaptive, capturing nuanced aspects of reasoning quality that are challenging to specify manually. Grounding the reward function in these structured rubrics allows RiT to guide the model toward producing trajectories that are simultaneously high-quality and interpretable, establishing a principled foundation for subsequent RL.

Response Reward. This reward evaluates the quality of the model’s final response and adapts to different task types:

- **Open-ended generation tasks:** For tasks without a unique correct answer, we employ **response rubrics**, analogous to the thinking rubrics, to assess multiple dimensions of the output. These rubrics are constructed according to the evaluation criteria of the underlying dataset, ensuring alignment with task-specific standards. A frontier LLM critic assigns binary scores for each rubric, and the aggregated response reward is computed as:

$$R_{\text{resp}} = \frac{1}{M} \sum_{m=1}^M s_m^{\text{resp}}, \quad (5)$$

where $s_m^{\text{resp}} \in \{0, 1\}$ denotes whether the response meets the m -th criterion of M rubrics.

- **Reasoning tasks:** For tasks with a well-defined answer, the response reward is based on **accuracy**, i.e., whether the model’s output matches the reference answer. Formally,

$$R_{\text{resp}} = \mathbf{1}\{o_{\text{model}} = o_{\text{gold}}\}, \quad (6)$$

where o_{model} and o_{gold} denote the model-generated and gold answer, respectively.

Gated Fusion Reward. To jointly optimize both the reasoning process and the final response, we integrate the thinking reward R_{thinking} and the response reward R_{resp} through a weighted fusion inspired by (Yu et al., 2025b; Peng et al., 2025), followed by a gated mechanism that enforces balanced performance. Specifically, let $\alpha \in [0, 1]$ denote the weight assigned to the reasoning component, with $1 - \alpha$ corresponding to the response component:

$$R_{\text{fusion}} = \alpha R_{\text{thinking}} + (1 - \alpha) R_{\text{resp}}. \quad (7)$$

The fused score is then passed through a **gating function** to obtain the final reward:

$$R_{\text{final}} = \text{Gate}(R_{\text{fusion}}, R_{\text{resp}}). \quad (8)$$

In this work, we instantiate the gating function as a hard minimum, such that the final reward is bounded by the weaker of the two components. The gating mechanism prevents degenerate optimization by adaptively constraining the influence of reasoning signals. This design ensures that high reward can only be achieved when both the reasoning

process and the final response exhibit strong performance. Consequently, the model is discouraged from over-optimizing either reasoning or output in isolation, while retaining flexibility for future extensions with soft or adaptive gating mechanisms.

4.4 Reinforcement Learning with GRPO

To optimize the proposed reward formulation, we adopt GRPO as the RL algorithm. Given a query q , the policy π_θ samples a group of candidate responses $\{y_i\}_{i=1}^g$. Each response is assigned a scalar reward $R(q, y_i)$ computed using the proposed gated fusion of thinking and response scores. GRPO then calculates a group-relative advantage by normalizing rewards within each group:

$$A_i = \frac{R(q, y_i) - \bar{R}_g}{\text{std}_{j \in g}(R(q, y_j))}, \quad (9)$$

where $\bar{R}_g = \frac{1}{|g|} \sum_{i \in g} R(q, y_i)$ and std denotes the standard deviation over the group. Based on the normalized advantage A_i , the policy is optimized by maximizing a clipped surrogate objective:

$$\mathcal{L}_{\text{GRPO}}(\theta) = \mathbb{E}_i \left[\min(\rho_i(\theta)A_i, \text{clip}(\rho_i(\theta), 1 - \epsilon, 1 + \epsilon)A_i) \right], \quad (10)$$

where $\rho_i(\theta) = \frac{\pi_\theta(y_i|q)}{\pi_{\theta_{\text{old}}}(y_i|q)}$, ϵ is a hyperparameter controlling the clipping range to prevent excessively large policy updates.

5 Experiments

In this section, we conduct a detailed evaluation of our method RiT. Implementation details are provided in the Appendix A.1 and the source code is attached as supplementary materials.

5.1 Experimental Setup

Datasets. We evaluate RiT on diverse datasets spanning open-ended generation and reasoning tasks. For open-ended generation, we use **LongBench-Write** (Bai et al., 2025), which targets long-form text requiring coherent and contextually rich outputs, and **WritingBench** (Wu et al., 2025), a multi-domain dataset covering six professional domains, including technical, creative, and business writing. For reasoning, we employ **StrategyQA** (Geva et al., 2021), consisting of multi-step strategic questions; **LogiQA** (Liu et al., 2023a), designed for logical deduction and inference; and **TruthfulQA** (Lin et al., 2022), which assesses the truthfulness of model-generated answers. RiT is

trained on the available training sets and evaluated on the corresponding test sets. Dataset statistics and details are provided in the Appendix A.2.

Backbone Reasoning Models. We conduct experiments using three backbone reasoning models: **Qwen3-8B** (Yang et al., 2025), **Qwen3-1.7B** (Yang et al., 2025), and **DeepSeek-R1-Distill-Llama-8B** (DeepSeek-AI, 2025). These models cover different parameter scales and architectures, providing a representative set for evaluating the effectiveness and generality of our thinking rubric supervision.

Baselines. To rigorously assess the effectiveness of RiT, we consider two representative baselines. First, we evaluate the **base model**, which serves as a performance lower bound. Second, we construct a **GRPO baseline with outcome-only reward (ORM)**, where optimization depends exclusively on response quality and provides no explicit signal on the thinking process. This baseline shares the same optimization algorithm as our method but differs only in reward design. Importantly, our goal is not to achieve state-of-the-art performance, but to validate the utility of the thinking rubric reward.

Evaluation Metrics. For reasoning tasks, we use the standard Pass@1 (Chen et al., 2021) metric. For open-ended generation, we employ Qwen3-235B-A22B-Instruct (Yang et al., 2025) as the judge for all benchmarks. More details of evaluation are provided in the Appendix A.3

5.2 Main Results

Table 1 reports the main experimental results across different reasoning models and benchmark settings. Overall, RiT consistently outperforms all baselines across diverse evaluation scenarios, demonstrating strong generalization. In open-ended generation settings, it exhibits improved handling of long-form composition, producing responses that better satisfy complex writing conditions. In reasoning-oriented tasks, it achieves higher accuracy and stronger factual reasoning capability, reflecting more reliable and consistent decision-making behavior. Importantly, these performance gains are observed across models with varying parameter scales and architectural designs, ranging from lightweight models to larger, more expressive architectures. This suggests that RiT is not tightly coupled to a specific model size or backbone, and can be effectively applied across different model families, highlighting its robustness and scalability.

Reasoning Models	Methods	Open-Ended Generation		Reasoning		
		LongBench-Write	WritingBench	StrategyQA	LogiQA	TruthfulQA
DeepSeek-R1-Distilled -Llama-8B	Base	54.9	4.07	0.563	0.463	0.504
	ORM	<u>59.5</u>	<u>4.46</u>	<u>0.659</u>	<u>0.558</u>	<u>0.579</u>
	RiT	60.3	4.70	0.703	0.564	0.695
Qwen3-1.7B	Base	74.5	6.37	0.638	0.571	0.500
	ORM	<u>76.3</u>	<u>7.10</u>	<u>0.716</u>	<u>0.597</u>	<u>0.848</u>
	RiT	79.3	7.13	0.734	0.606	0.854
Qwen3-8B	Base	84.4	7.45	0.729	0.594	0.744
	ORM	<u>85.7</u>	<u>8.10</u>	<u>0.808</u>	<u>0.690</u>	<u>0.945</u>
	RiT	89.2	8.24	0.838	0.719	0.951

Table 1: Comparative evaluation of different methods across reasoning models and datasets. Best ones among all methods are **bold-faced**, and second-best are underlined.

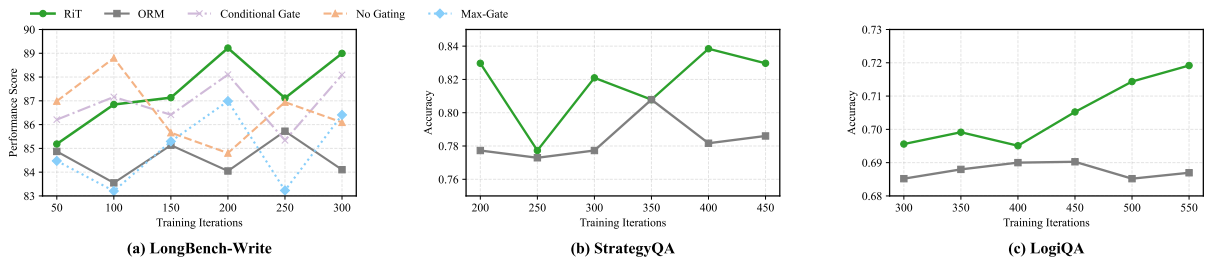


Figure 3: (a) LongBench-Write: ablation on both the thinking-rubric reward and the choice of gating function. (b) StrategyQA and (c) LogiQA: ablation on the thinking-rubric reward.

Notably, ORM can be viewed as an ablation of RiT without the thinking-rubric reward. Comparing this variant with RiT underscores the contribution of the thinking-rubric guidance, demonstrating that incorporating reasoning-quality feedback significantly improves performance across diverse tasks. We further observe that models with stronger base capabilities benefit more from the thinking-rubric guidance. This is likely because high-capacity models can better leverage the fine-grained feedback to produce more coherent, accurate, and complex outputs, whereas lower-capacity models may struggle to fully exploit the reward signal.

5.3 Ablation Study

Module Analysis. Table 2 presents an ablation study on Qwen3-8B across three benchmarks: LongBench-Write, StrategyQA, and LogiQA. Our analysis reveals four main insights: (1) The full RiT consistently outperforms all ablation variants, demonstrating that TR (Thinking-Rubric Reward) and RR (Response Reward) are not redundant but mutually reinforcing. (2) RR acts as the primary driver of performance in reasoning tasks, highlighting that direct feedback on the final answer is es-

essential for ensuring logical correctness and correcting terminal errors. (3) For open-ended tasks that lack a unique "gold" label, TR provides crucial guidance, effectively serving as a "compass" for intermediate reasoning. (4) Replacing GRPO with standard PPO (Schulman et al., 2017) results in a performance drop, suggesting that GRPO is better suited to balancing the heterogeneous reward scales introduced by the combination of thinking and response-level supervisions.

Fine-grained Ablation of the Thinking-Rubric Reward. To further analyze the effect of the thinking-rubric reward, we evaluate Qwen3-8B on three representative benchmarks, including LongBench-Write, StrategyQA, and LogiQA. We compare RiT with ORM using checkpoints near reward convergence, enabling a controlled examination of their performance trajectories. As observed in the Figure 3, during the late convergence stages, RiT maintains performance scores that are generally not lower than ORM, with relatively minor fluctuations. This indicates that RiT demonstrates relatively stable superiority during the convergence phase, suggesting both high performance and more reliable behavior as training stabilizes.

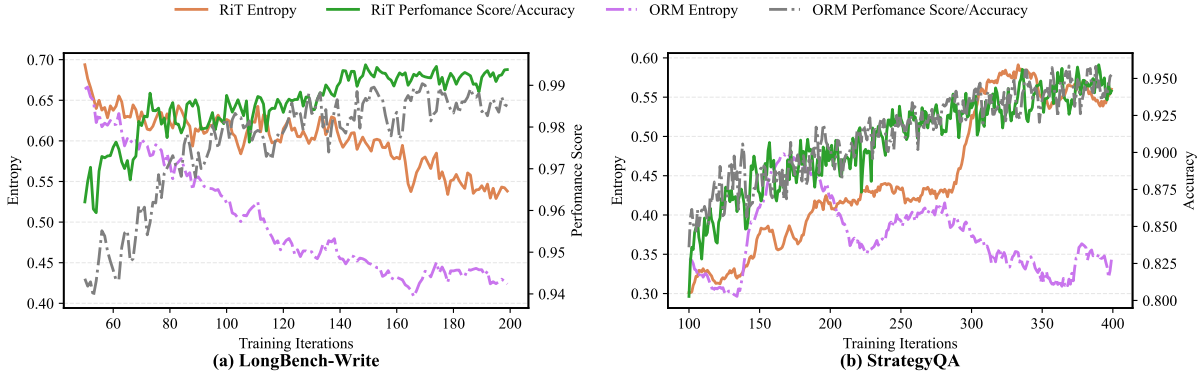


Figure 4: Entropy and performance of Qwen3-8B with RiT and ORM on LongBench-Write (a) and StrategyQA (b).

Impact of Gate Function Choice. We further analyze the impact of different gating strategies through ablation experiments on Qwen3-8B evaluated on the LongBench-Write, as illustrated in Figure 3 (a). Specifically, we compare several alternatives: (1) **No Gating**, where the final reward is simply the weighted fusion of thinking rubric and response scores; (2) a **Max-Gate**, which selects the larger of the two scores; and (3) a **Conditional Gate**, where the thinking-rubric reward is incorporated only when the response is fully correct (i.e., the response score equals 1). Empirically, we observe that most gated variants outperform the outcome-only baseline, indicating that incorporating thinking rubric signals generally benefits learning. However, among all designs, the Min-Gate (RiT) consistently yields the most stable and reliable improvements. In contrast, the Max-Gate often leads to degenerate behaviors, as high thinking rubric scores can dominate even when the final answer is incorrect or low quality, weakening coupling between thinking quality and task success.

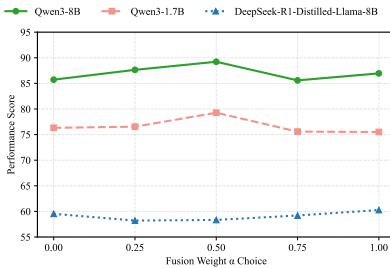


Figure 5: Performance scores of different models at varying fusion weights α in LongBench-Write.

Impact of Fusion Weight α Choice. Figure 5 illustrates the performance scores of different models under varying fusion weights α in LongBench-Write. Models with relatively higher baseline per-

Methods	LongBench-Write	StrategyQA	LogiQA
Base	84.4	0.729	0.594
RiT (w/o TR)	85.7	<u>0.808</u>	<u>0.690</u>
RiT (w/o RR)	85.8	0.769	0.638
RiT (PPO)	<u>88.0</u>	0.799	0.638
RiT	89.2	0.838	0.719

Table 2: Ablation on Qwen3-8B across LongBench-Write, StrategyQA, and LogiQA. TR and RR denote the Thinking-Rubric and Response rewards, respectively. Best results are in **bold**, and second-best are underlined.

formance generally tend to achieve their best results at moderate α values, as excessive weighting toward the thinking score can override the model’s strong predictions and slightly degrade performance. In contrast, the model with lower baseline performance, such as DeepSeek-R1-Distill-Llama-8B, benefits from higher α , suggesting that the thinking score can compensate for weaker predictions and consequently enhance overall results.

5.4 Entropy Analysis

We further analyze the entropy dynamics of Qwen3-8B on the LongBench-Write and StrategyQA datasets, with a particular focus on the near-convergence stage of training. As illustrated in Figure 4, on LongBench-Write, the entropy of both methods exhibits an overall decreasing trend as training progresses. Nevertheless, RiT demonstrates a noticeably slower rate of entropy reduction compared to ORM, indicating its ability to preserve a higher degree of exploration during optimization. In contrast, on the reasoning-oriented StrategyQA task, RiT maintains an increasing entropy trend even near convergence, whereas ORM exhibits a clear decline in entropy. This suggests that RiT is able to sustain higher uncertainty and flexibility in its reasoning process, which may facilitate

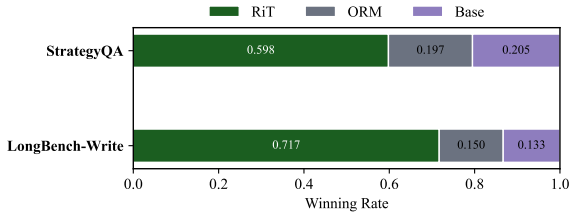


Figure 6: Comparison of the thinking quality of Qwen3-8B under three strategies (Base, ORM, and RiT) on LongBench-Write and StrategyQA. A frontier LLM selects the best reasoning trajectory for each query, and the fraction of selections (i.e., the winning rate) is used as a proxy for the thinking quality.

the exploration of diverse solution paths. Overall, these observations suggest that incorporating the thinking rubric endows the model with stronger exploratory capacity, enabling more effective reasoning behaviors across different task types. This enhanced ability to explore and deliberate allows the model to achieve improved performance in both open-ended generation and reasoning benchmarks.

5.5 Thinking Quality Comparison

We evaluate the thinking quality of Qwen3-8B under three strategies: Base, ORM, and RiT on LongBench-Write and StrategyQA. Following prior work, a frontier LLM (i.e., Claude Sonnet 4 (Anthropic, 2025)) selects the best reasoning trajectory for each query, and the selection ratio is used as the thinking quality metric. As shown in Figure 6, RiT consistently achieves the highest selection rate across both datasets, indicating more coherent and goal-aligned reasoning. This trend aligns with our entropy analysis, suggesting that thinking-guided training enhances exploratory capacity and overall reasoning quality beyond outcome-only supervision.

5.6 Case Study

We conduct case studies to qualitatively assess the impact of thinking-rubric supervision, with a representative example from LongBench-Write shown in Figure 7 (additional cases are provided in Appendix A.4). In a task requiring a 1,200-word generation, ORM fails to regulate content expansion, resulting in a 2,119-word output, which substantially exceeds the target length. In contrast, RiT produces a more aligned 1,328-word response, closer to the required constraint. This precision stems directly from the length-planning dimension within our thinking rubrics: RiT’s thinking process ex-

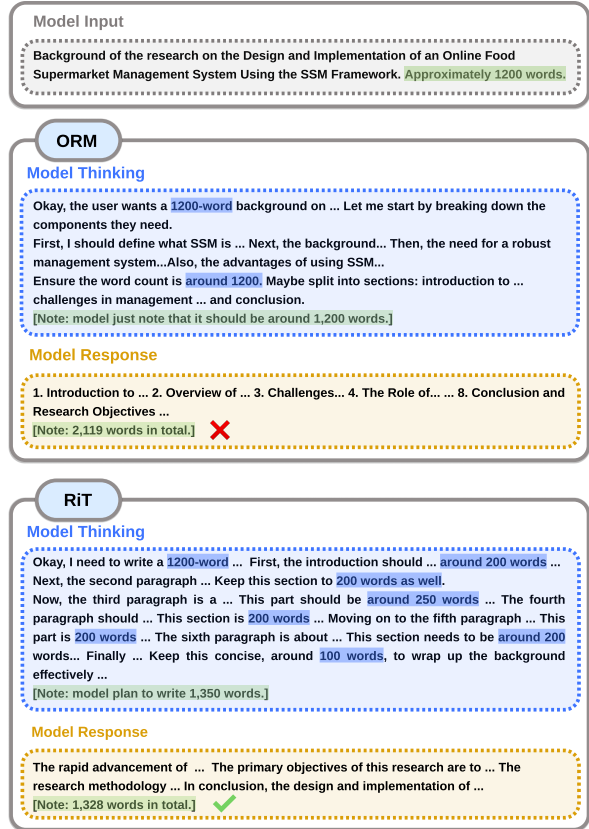


Figure 7: A case study on open-ended generation showing that ORM exceeds the target length due to coarse planning, while RiT meets the length requirement by explicitly planning word budgets during thinking.

Explicitly pre-allocates word budgets to each subsection. This demonstrates that under thinking rubrics guided supervision of intermediate planning, RiT fosters the disciplined, constraint-aware reasoning necessary for complex, open-ended tasks.

6 Conclusion

In this paper, we propose RiT, the first framework to incorporate rubric-based supervision into intermediate reasoning, enabling direct evaluation and optimization of thinking trajectories. RiT automatically generates fine-grained thinking rubrics using and assigns thinking-level rewards based on rubric compliance. To balance reasoning quality with final-task performance, we integrate the thinking reward with an outcome-based reward through a gated fusion mechanism, ensuring consistency between intermediate reasoning and final responses. Experiments demonstrate that RiT improves reasoning quality across diverse tasks, highlighting the value of thinking-rubrics as a scalable and flexible framework for guiding human-aligned reasoning.

560 Limitations

561 Despite the effectiveness of RiT, several limitations
562 remain. First, the generation of thinking rubrics
563 relies automated processes with minimal human
564 oversight, which may limit their alignment with
565 nuanced human expectations or domain-specific
566 knowledge. Incorporating human-in-the-loop ver-
567 ification could help mitigate this issue in future
568 work. Second, the integration of thinking-rubric
569 and response rewards currently requires manual
570 hyperparameter tuning. Future work could explore
571 adaptive mechanisms to dynamically balance these
572 reward scales across tasks of varying complexity.
573 Finally, due to computational constraints, experi-
574 ments were conducted on moderately sized models
575 (e.g., 8B parameters), leaving the scalability and
576 emergent behaviors of RiT on ultra-large models
577 (e.g., >70B parameters) to be investigated.

578 References

579 Anthropic. 2025. Introducing claude 4. [https://www.](https://www.anthropic.com/news/claude-4)
580 [anthropic.com/news/claude-4](https://www.anthropic.com/news/claude-4).

581 Ge Bai, Jie Liu, Xingyuan Bu, Yancheng He, Jia-
582 heng Liu, Zhanhui Zhou, Zhuoran Lin, Wenbo Su,
583 Tiezheng Ge, Bo Zheng, and Wanli Ouyang. 2024.
584 [Mt-bench-101: A fine-grained benchmark for evalu-](#)
585 [ating large language models in multi-turn dialogues.](#)
586 *In Proceedings of the 62nd Annual Meeting of the*
587 *Association for Computational Linguistics (Volume*
588 *1: Long Papers), ACL 2024, Bangkok, Thailand, Au-*
589 *gust 11-16, 2024*, pages 7421–7454. Association for
590 Computational Linguistics.

591 Yushi Bai, Jiajie Zhang, Xin Lv, Linzhi Zheng, Siqi
592 Zhu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li.
593 2025. [Longwriter: Unleashing 10,000+ word gener-](#)
594 [ation from long context llms.](#) *In The Thirteenth In-*
595 *ternational Conference on Learning Representations,*
596 *ICLR 2025, Singapore, April 24-28, 2025.* OpenRe-
597 view.net.

598 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,
599 Henrique Pondé de Oliveira Pinto, Jared Kaplan,
600 Harri Edwards, Yuri Burda, Nicholas Joseph, Greg
601 Brockman, Alex Ray, Raul Puri, Gretchen Krueger,
602 Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela
603 Mishkin, Brooke Chan, Scott Gray, and 39 others.
604 2021. [Evaluating large language models trained on](#)
605 [code.](#) *CoRR*, abs/2107.03374.

606 Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng,
607 Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang
608 Zhou, Te Gao, and Wanxiang Che. 2025. [To-](#)
609 [wards reasoning era: A survey of long chain-of-](#)
610 [thought for reasoning large language models.](#) *CoRR*,
611 abs/2503.09567.

DeepSeek-AI. 2025. [Deepseek-r1: Incentivizing rea-](#)
612 [soning capability in llms via reinforcement learning.](#)
613 *CoRR*, abs/2501.12948. 614

Ahmed El-Kishky, Alexander Wei, Andre Saraiva, Bo-
615 rys Minaiev, Daniel Selsam, David Dohan, Francis
616 Song, Hunter Lightman, Ignasi Clavera Gilaberte,
617 Jakub Pachocki, Jerry Tworek, Lorenz Kuhn, Lukasz
618 Kaiser, Mark Chen, Max Schwarzer, Mostafa Roha-
619 ninejad, Nat McAleese, o3 contributors, Oleg Mürk,
620 and 5 others. 2025. [Competitive programming with](#)
621 [large reasoning models.](#) *CoRR*, abs/2502.06807. 622

Lishui Fan, Yu Zhang, Mouxiang Chen, and Zhongxin
623 Liu. 2025. [Posterior-grpo: Rewarding reasoning pro-](#)
624 [cesses in code generation.](#) *CoRR*, abs/2508.05170. 625

Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot,
626 Dan Roth, and Jonathan Berant. 2021. [Did aristotle](#)
627 [use a laptop? A question answering benchmark with](#)
628 [implicit reasoning strategies.](#) *Trans. Assoc. Comput.*
629 *Linguistics*, 9:346–361. 630

Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan,
631 Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen,
632 Shengjie Ma, Honghao Liu, Yuanzhuo Wang, and
633 Jian Guo. 2024. [A survey on llm-as-a-judge.](#) *CoRR*,
634 abs/2411.15594. 635

Anisha Gunjal, Anthony Wang, Elaine Lau, Vaskar
636 Nath, Bing Liu, and Sean Hendryx. 2025. [Rubrics as](#)
637 [rewards: Reinforcement learning beyond verifiable](#)
638 [domains.](#) *CoRR*, abs/2507.17746. 639

Zenan Huang, Yihong Zhuang, Guoshan Lu, Zeyu Qin,
640 Haokai Xu, Tianyu Zhao, Ru Peng, Jiaqi Hu, Zhan-
641 ming Shen, Xiaomeng Hu, Xijun Gu, Peiyi Tu, Jiaxin
642 Liu, Wenyu Chen, Yuzhuo Fu, Zhiting Fan, Yanmei
643 Gu, Yuanyuan Wang, Zhengkai Yang, and 2 others.
644 2025. [Reinforcement learning with rubric anchors.](#)
645 *CoRR*, abs/2508.12790. 646

Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richard-
647 son, Ahmed El-Kishky, Aiden Low, Alec Hel-
648 yar, Aleksander Madry, Alex Beutel, Alex Carney,
649 Alex Ifitimie, Alex Karpenko, Alex Tachard Passos,
650 Alexander Neitz, Alexander Prokofiev, Alexander
651 Wei, Allison Tam, Ally Bennett, Anyanya Kumar, and
652 80 others. 2024. [Openai o1 system card.](#) *CoRR*,
653 abs/2412.16720. 654

Ruipeng Jia, Yunyi Yang, Yongbo Gai, Kai Luo, Shi-
655 hao Huang, Jianhe Lin, Xiaoxi Jiang, and Guan-
656 jun Jiang. 2025. [Writing-zero: Bridge the gap between](#)
657 [non-verifiable tasks and verifiable rewards.](#) *CoRR*,
658 abs/2506.00103. 659

Muhammad Khalifa, Rishabh Agarwal, Lajanugen Lo-
660 geswaran, Jaekyeom Kim, Hao Peng, Moontae Lee,
661 Honglak Lee, and Lu Wang. 2025. [Process reward](#)
662 [models that think.](#) *CoRR*, abs/2504.16828. 663

Jinu Lee and Julia Hockenmaier. 2025. [Evaluating](#)
664 [step-by-step reasoning traces: A survey.](#) *CoRR*,
665 abs/2502.12289. 666

780 Mengyue Wu, Qin Jin, and Fei Huang. 2025. [Writingbench: A comprehensive benchmark for generative writing](#). *CoRR*, abs/2503.05244.

781

782

783 Fengli Xu, Qianyu Hao, Zefang Zong, Jingwei Wang,
784 Yunke Zhang, Jingyi Wang, Xiaochong Lan, Jiahui
785 Gong, Tianjian Ouyang, Fanjin Meng, Chenyang
786 Shao, Yuwei Yan, Qinglong Yang, Yiwen Song, Si-
787 jian Ren, Xinyuan Hu, Yu Li, Jie Feng, Chen Gao,
788 and Yong Li. 2025. [Towards large reasoning models: A survey of reinforced reasoning with large language models](#). *CoRR*, abs/2501.09686.

789

790

791 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,
792 Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,
793 Chengen Huang, Chenxu Lv, Chujie Zheng, Day-
794 iheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao
795 Ge, Haoran Wei, Huan Lin, Jialong Tang, and 40
796 others. 2025. [Qwen3 technical report](#). *CoRR*,
797 abs/2505.09388.

798

799 Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan,
800 Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu,
801 Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole
802 Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang,
803 Mofan Zhang, Wang Zhang, Hang Zhu, and 16 others.
804 2025a. [DAPO: an open-source LLM reinforcement learning system at scale](#). *CoRR*, abs/2503.14476.

805

806 Zhuohao Yu, Weizheng Gu, Yidong Wang, Xingru Jiang,
807 Zhengran Zeng, Jindong Wang, Wei Ye, and Shikun
808 Zhang. 2025b. [Reasoning through execution: Unifying process and outcome rewards for code generation](#). In *Forty-second International Conference on Machine Learning, ICML 2025, Vancouver, BC, Canada, July 13-19, 2025*. OpenReview.net.

809

810

811

812 Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui
813 Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong
814 Liu, Rui Men, An Yang, Jingren Zhou, and Jun-
815 yang Lin. 2025. [Group sequence policy optimization](#).
816 *CoRR*, abs/2507.18071.

817

818 Yang Zhou, Sunzhu Li, Shunyu Liu, Wenkai Fang, Jiale
819 Zhao, Jingwen Yang, Jianwei Lv, Kongcheng Zhang,
820 Yihe Zhou, Hengtong Lu, Wei Chen, Yan Xie, and
821 Mingli Song. 2025. [Breaking the exploration bottleneck: Rubric-scaffolded reinforcement learning for general LLM reasoning](#). *CoRR*, abs/2508.16949.

822

A Details of Experiments

A.1 Implementation Details

We train our models using the VeRL framework¹ (Sheng et al., 2025) and follow standard GRPO hyperparameter settings. Specifically, the learning rate is set to 2×10^{-6} , and each rollout consists of 8 sampled responses. For open-ended generative tasks, we use a batch size of 32, while for reasoning or inference tasks, the batch size is set to 128. The thinking rubrics are generated by Claude Sonnet 4 (Anthropic, 2025). During training, model outputs are scored by Qwen3-235B-A22B-Instruct. For open-ended generation, the fusion weight α is generally set to 0.5, while for DeepSeek-R1-Distilled-Llama-8B on LongBench-Write, α is set to 1. For reasoning tasks, α is consistently set to 1. Training proceeds until the reward converges. The models are run for around 500 steps on reasoning tasks, and around 200 steps for open-ended generation tasks.

A.2 Datasets

Here we detail the training and test splits for all datasets used in our experiments. Different split ratios (50:50 vs. 80:20) are chosen based on dataset size and evaluation needs, while ensuring that training and test data are drawn from the same underlying distribution. An overview of the splits is provided in Table 3.

- **LongBench-Write** is split 50:50 by type and text length, preserving the original distribution across these attributes to ensure balanced evaluation of long-form generation quality.
- **WritingBench** is split 50:50 via stratified sampling over major/minor domains and language, ensuring representative coverage of all six professional writing domains in both splits.
- **StrategyQA** uses the official benchmark split, with the training set and development set as the test set for comparability with prior work.
- **LogiQA** adopts an NLI task formulation as suggested by (Liu et al., 2023a), focusing on logical inference between premise and hypothesis. We use the benchmark-provided training and test sets to ensure consistency with the standard evaluation protocol.

- **TruthfulQA** uses the multiple-choice 1 (MC1) task format and is split 80:20 into training and test sets. This ratio provides sufficient training data given the limited dataset size while preserving a reliable test set.

Dataset	Training Set	Test Set
LongBench-Write	60	60
WritingBench	500	500
StrategyQA	2,061	229
LogiQA	30,908	3,942
TruthfulQA	653	164

Table 3: Overview of the training and test splits for all datasets used in our experiments.

A.3 Evaluation Protocol

Here we describe the evaluation protocol and scoring procedures used in our experiments.

For LongBench-Write and WritingBench, we use the official evaluation protocols and directly adopt the evaluation scripts released in their public GitHub repositories², without any modification to the evaluation prompts or scoring logic.

For LongBench-Write, the final score ranges from 0 to 100 and is computed as the average of a length score (S_l) and a quality score (S_q). For WritingBench, the evaluator independently assigns a 10-point score to the generated response for each quality criterion $c_i \in C_q$, and the final score is computed by averaging across all criteria.

A.4 Case Study

We present additional case studies to illustrate the effectiveness of thinking-rubric supervision in guiding model reasoning and content generation.

WritingBench (Open-Ended Generation). As shown in Figure 8, the task is to write a detailed guidance document for parents, incorporating interview summaries and industry reports. ORM lacks explicit planning for integrating external inputs, resulting in omissions and disorganized content. RiT leverages content-mapping from the thinking rubrics, specifying which information belongs in which section, and produces a coherent, information-complete document.

²<https://github.com/THUDM/LongWriter> and <https://github.com/X-PLUG/WritingBench>

¹<https://github.com/volcengine/verl>

Model Input
Write a detailed guidance document on how parents can work with educational consulting agencies to develop personalized learning plans for their children. The document should include but not be limited to: ... Interview summaries from experts or educational consultants, or public industry reports Summary of Interviews with Experts or Education Consultants and Public Industry Reports: ...

ORM

Model Thinking
... The user also provided interview summaries and industry reports. I need to integrate that information into the document. For example, the future of the education consulting industry, market growth points, integration with the internet, and talent training. Maybe add a section on industry trends and how they affect the approach parents take. ...
[Note: No clear plan for placement]

Model Response
Introduction: The Role of Educational Consulting in Personalized Learning ... Section 1: Preparing for Engagement with Educational Consulting Agencies ... Section 2: Initial Communication with Educational Consulting Agencies ... Section 3: The Process of Developing a Personalized Learning Plan ...
[Note: Missing summaries of interviews and industry reports] ❌

RiT

Model Thinking
... The industry reports they provided are about the future of education consulting in China, market growth, integration with the internet, talent development, etc. I need to weave these into the document where relevant, especially in sections discussing the industry's growth and the role of consulting agencies. ... The industry reports mention the growth of the education consulting industry, integration with the internet, and talent development. I need to connect these points to the document. For instance, when discussing the future of the industry, mention how technology integration (like "Internet+" platforms) can enhance personalized learning. Also, when talking about talent, emphasize the importance of professional consultants and the need for ongoing training ...
[Note: Specifies content and maps it to sections]

Model Response
VII. The Role of the Education Consulting Industry in Future Development... Integration with Technology ("Internet+")... Segmented Services... Talent Development... System Integration... Example: ...
[Note: Interviews and industry reports have been integrated and placed in Section VII as planned.] ✅

Figure 8: A case study on open-ended generation. ORM lacks structured planning and omits key details, while RiT uses content-mapping from thinking rubrics to map content to sections and produce a coherent document.

Model Input
Are all the elements plants need for photosynthesis present in atmosphere of Mars?

ORM

Model Thinking
Okay, let's tackle this question ...
Wait, the question is about elements needed for photosynthesis. Let me think about the elements. Carbon, hydrogen, oxygen ...
Wait, but the question is about elements plants need. Let me think about the elements ...
So, if the elements are present, even if in low amounts, the answer is yes. Therefore, the answer is yes.
[Note: Repeated thinking]

Model Response
Answer: No.
[Note: The thinking and the response are inconsistent] ❌

RiT

Model Thinking
Okay, so the question is whether all the elements plants need for photosynthesis are present in Mars' atmosphere. Let me break this down step by step ...
But I'm still uncertain because I remember that Mars' atmosphere has very little water, but the question is about presence.
...
Thus, the answer is Yes.

Model Response
Answer: Yes.
[Note: The thinking and the response are consistent] ✅

Figure 9: A case study on a reasoning task. ORM exhibits repetitive and unfocused reasoning, producing an incorrect answer, whereas RiT follows a disciplined reasoning process guided by thinking rubrics and arrives at the correct answer.

StrategyQA (Reasoning Task). As shown in Figure 9, the task asks whether Mars' atmosphere contains all elements needed for photosynthesis. ORM exhibits looping reasoning, repeatedly revisiting similar content and ultimately producing an incorrect answer despite internally concluding correctly. RiT, guided by thinking-rubric supervision, maintains a disciplined reasoning trajectory, prioritizing critical elements and avoiding redundant loops, leading to the correct answer.

Overall, these case studies show thinking-rubric supervision enables RiT to produce more precise, structured, and constraint-aware outputs across reasoning and open-ended generation tasks.