

RAVR: Reference-Answer-Guided Variational Reasoning for Large Language Models

Tianqianjin Lin^{1*}, Xi Zhao², Xingyao Zhang², Rujiao Long², Yi Xu²,
Zhuoren Jiang¹, Wenbo Su², Bo Zheng^{2†}

¹Zhejiang University, ²Alibaba Group
lintqj@zju.edu.cn

Abstract

Reinforcement learning (RL) can refine the reasoning abilities of large language models (LLMs), but critically depends on a key prerequisite: the LLM can already generate high-utility reasoning paths with non-negligible probability. For tasks beyond the LLM’s current competence, such reasoning paths can be hard to sample, and learning risks reinforcing familiar but suboptimal reasoning. We are motivated by the insight from cognitive science that *Why is this the answer* is often an easier question than *What is the answer*, as it avoids the heavy cognitive load of open-ended exploration, opting instead for explanatory reconstruction—systematically retracing the reasoning that links a question to its answer. We show that LLMs can similarly leverage answers to derive high-quality reasoning paths. We formalize this phenomenon and prove that conditioning on the answer provably increases the expected utility of sampled reasoning paths, thereby transforming intractable problems into learnable ones. Building on this insight, we introduce **RAVR**¹ (Reference-Answer-guided Variational Reasoning), an end-to-end framework that uses answer-conditioned reasoning as a variational surrogate for question-only reasoning. Experiments across 12 benchmarks and 4 models demonstrate the effectiveness of **RAVR**, and analysis of the reasoning behavior shows that **RAVR** reduces hesitation, strengthens conclusion consolidation, and promotes problem-specific strategies in reasoning.

1 Introduction

Large language models (LLMs) can solve increasingly complex problems when guided by reinforcement learning (RL) (Zhang et al., 2025a). In this

realm, a trajectory is a completion consisting of a reasoning path followed by a final response (OpenAI, 2024; DeepSeek-AI, 2025; Qwen, 2025; OpenAI, 2025). The objective is straightforward: to sample completions from the model’s distribution and then shift probability mass toward those with higher advantage, such as ones that produce a correct answer. This process is more like redistributing probabilities among sampleable completions, rather than generating new ones (Yue et al., 2025). This indicates a critical prerequisite for effective optimization—the model must already be able to sample useful completions with non-negligible probability. In popular relative-advantage approaches such as GRPO (Shao et al., 2024), this prerequisite becomes even stricter, since the advantage of each completion is defined relative to others—meaning that even weak completions can be reinforced as long as they are better than the rest. Unfortunately, tasks beyond the model’s competence or outside its preferences make high-utility completions difficult to obtain (Zhang et al., 2025b; Li et al., 2025). As a result, training collapses into reinforcing a narrow set of familiar but suboptimal completions, while promising ones remain unexplored.

To address this issue, we advance a simple thesis: utilizing the reference answer helps derive good reasoning paths. While the reference answer is available in the training data, current methods use it only to compute reward; we argue that its potential can be more fully exploited. In cognitive science, *Why is this the answer* is often an easier question than *what is the answer* because it relieves the learner from high cognitive load of *open-ended exploration* based solely on the question. Instead, it allows the learner to concentrate on *explanatory reconstruction*—tracing the logic that connects the question to the reference answer (Chi et al., 1989). For example, access to the answer can help learners detect errors to keep exploration, avoid overthinking when they are correct, and even engage

*This work was done during the author’s internship at the Future Life Lab, Alibaba Group.

†Corresponding Author.

¹The code is available at <https://github.com/Lintianqianjin/RAVR>.

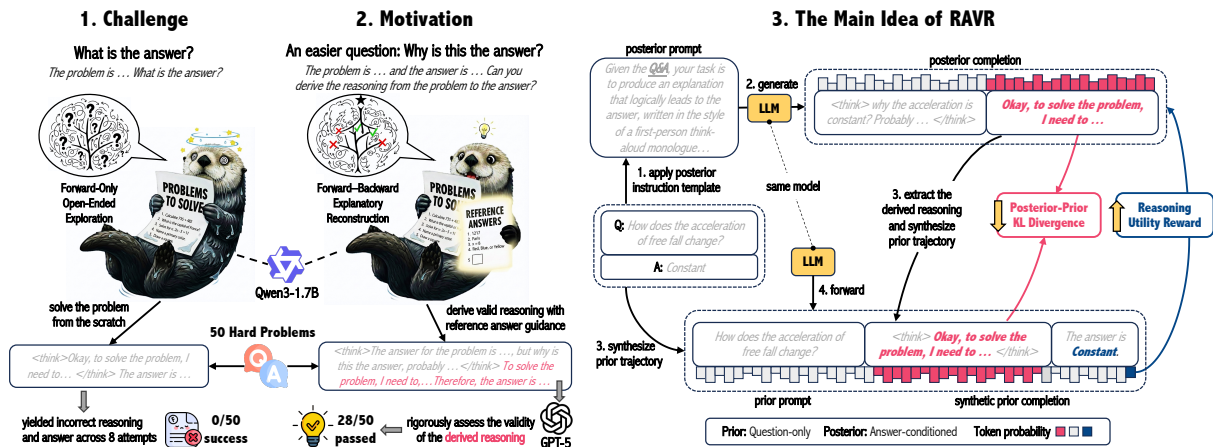


Figure 1: We sampled 50 hard questions from CrossThink-QA (Akter et al., 2026) on which Qwen3-1.7B, with thinking mode enabled, failed to obtain the correct answer across 8 attempts. For each question, we provided the answer in the prompt and asked the LLM to derive the reasoning. A GPT-5 judge, with thinking mode enabled, using strict criteria, evaluated whether the derived reasoning was logically coherent without indicating access to answer. In over 50% of the unlearnable questions, the LLM was able to derive valid reasoning. See Appendix A.2 for detailed prompts and cases. Based on this observation, we propose **RAVR**. We use answer-conditioned posterior to help the learning of the question-only prior. The LLM is instructed to derive a reasoning path from the question to the answer with thinking mode enabled. **RAVR** regards the reference-answer token probability as the reward for the generated reasoning and maximize it to enhance the ability of the LLM to think *why is this the answer*. Meanwhile, **RAVR** minimizes the KL divergence between the posterior and the question-only prior to help the model better think *what is the answer* and in turn, the prior also regularizes the behavior of the posterior. See Section 2 for details.

in backward reasoning to obtain the preconditions leading to the answer. The LLM is likely capable of simulating this human behavior. We validate this intuition with a motivation experiment as shown in Figure 1: for many hard questions where the LLM fails to sample a valid reasoning to solve the problem correctly after multiple attempts, providing the answer enables it to generate rational reasoning.

Building on this insight, we propose using reference answer as a guide to better explore high-quality reasoning. We formalize this effect and show that conditioning on the answer provably increases the likelihood of sampling sound reasoning paths and reduces the likelihood of sampling flawed reasoning paths, thereby improving the overall expected quality of reasoning over the sampling distribution.

While this is conceptually clean, this reasoning generation pipeline does not align with the real need at inference—the LLM can only observe the question and needs to generate the reasoning instantly in the thinking tags. We address this issue from the perspective of variational inference (Gershman and Goodman, 2014), viewing the *derived reasoning* conditioned on the reference answer as a variational surrogate for the reasoning inside the thinking tags when question-only.

We optimize the evidence lower bound (Kingma and Welling, 2013) by maximizing expected reasoning utility over the answer-conditioned distribution while reducing the discrepancy between the answer-conditioned reasoning distribution (hereafter *Posterior*) and the question-only reasoning distribution (hereafter *prior*) via Kullback–Leibler (KL) divergence (derivation in Section 2.3).

To further stabilize training, we introduce the following designs. *First*, we leverage a role-play style prompt (Shanahan et al., 2023), asking the LLM to produce a first-person think-aloud monologue as if solving the problem from scratch. The monologue is allowed to include any classical reasoning behavior such as reflection and backtracking. We introduce this design to bridge the language style gap between the two distributions: while the *derived reasoning* generated after the thinking tags can be concise and direct, the internal reasoning under prior is usually exploratory and reflective. This setup encourages the *derived reasoning* to mirror the behavior expected inside the thinking tags. *Second*, we introduce a utility baseline for estimating reasoning utility under the posterior, which is the expected reasoning utility under the prior. This provides a more informative reward signal by measuring how much the reasoning from the

posterior improves over that from the prior. *Third*, we reweight samples in the KL estimation using the utility reward, aiming to align the prior to a distribution whose expected utility is as large as possible.

With the proposed variational objective and the designed strategies, we can end-to-end train a single LLM that learns to reason effectively from scratch at inference. We name this framework as **RAVR**, for **R**eference-**A**nswer-guided **V**ariational **R**easoning. Results on standard benchmarks for both general and math reasoning show that **RAVR** substantially enhances reasoning capabilities and outperforms state-of-the-art methods. For example, when we use **RAVR** to train Qwen3-1.7B on the CrossThink-QA (Akter et al., 2026) dataset, it achieves a GPQA-Diamond score of 40.91, outperforming DAPO (Yu et al., 2025a) by 5.56 points. Additionally, we analyze the reasoning behavior of the LLM trained with **RAVR** and find that **RAVR** can reduce hesitation, strengthen conclusion consolidation, and promote problem-specific strategies in reasoning.

Overall, this paper makes the following contributions: (1) We formalize the intuition that reference answers can guide reasoning. We prove that conditioning on the answer provably amplifies the probability of high-utility reasoning paths; (2) We propose **RAVR**, the first end-to-end framework that operationalizes this insight and also the first to leverage the reasoning ability of the LLM to use the reference answer. **RAVR** alleviates the exploration difficulty in RL for LLM; (3) We demonstrate the effectiveness of **RAVR** through extensive experiments on both general and math domains and the in-depth analysis of the reasoning behavior of the LLM trained with **RAVR**. We share our code to facilitate future research.

2 RAVR: Reference-Answer-Guided Variational Reasoning

In this section, we introduce the background of RL for LLM, explain the motivation of **RAVR** and then describe the proposed objective and strategies to realise the motivation.

2.1 Preliminaries and the Sampling Challenge

In this work, we consider training an LLM, denoted by π_θ , to be a Large Reasoning Model (LRM). Given a sample (x, y^*) , the LLM π_θ is supposed to produce an intermediate *reasoning path* z before

giving a final answer y for each problem x :

$$\pi_\theta(z, y | x) = \pi_\theta(z | x) \pi_\theta(y | x, z). \quad (1)$$

In general, a reward function provides a reward $R(y) \in \mathbb{R}$ to judge the consistency of the reference answer y^* and the generated answer y . The standard objective is to maximize the expected reward:

$$\mathcal{J}(\theta) = \mathbb{E}_{z \sim \pi_\theta(\cdot | x)} \mathbb{E}_{y \sim \pi_\theta(\cdot | x, z)} [R(y)]. \quad (2)$$

To align with the comparative nature of reward models, recent popular RL algorithms for LLM usually normalize the reward of each completion in a group-relative or batch-relative manner to obtain the advantage of each completion. Taking GRPO (Shao et al., 2024) as an example, it first generates multiple completions for a problem x , and calculates the advantage of i -th completion as $A_i = \frac{R_i - \text{mean}(R)}{\text{std}(R)}$, where mean and std represent the average and standard deviation of the rewards. *This indicates a critical prerequisite for effective optimization—the model must already be able to sample good completions with non-negligible probability*, because even weak completions can be reinforced as long as they are better than the rest and obtain a positive advantage. When few good completions are sampled, training risks collapsing into reinforcing familiar but suboptimal completions. If $R(\cdot)$ is defined as binarized correctness, which is common in current applications, this challenge can be more severe when no correct completion is generated since all the advantages become zero, making no optimization signal.

2.2 Motivation: Conditioning on Reference Answer Amplifies Good Reasoning

The reasoning path z dominates the final quality of the entire completion since it causally influences the generation of the final answer y . Therefore, we propose to mitigate the sampling challenge of high-utility reasoning paths to alleviate the sampling challenge of the overall completion. In this section, we formally define the utility of a reasoning path, and prove our motivation that conditioning on reference answer can amplify the sampling probability of high-utility reasoning paths.

If $R(\cdot)$ denotes binarized correctness, optimizing reasoning path z aims to maximize the expectation it yields the reference answer. Hence, a natural alternative is to take the LLM’s likelihood of the reference-answer as reward on z and the LLM can be optimized by maximizing the expected reward:

$$\mathcal{J}_{\text{prob}}(\theta) = \mathbb{E}_{z \sim \pi_{\theta}(\cdot|x)} [\pi_{\theta}(y^* | x, z)]. \quad (3)$$

This objective does not fully match the original setting, for example, the reference answer y^* is usually a single word or phrase such as an option ‘‘A’’ of a multiple-choice question, but the prediction y can be a paragraph that includes steps to reach the answer. Yet, recent works have shown its acceptable effectiveness (Yu et al., 2025b; Zhou et al., 2025b), especially in open-ended domain (Xu et al., 2025). Therefore, it provides a reliable operational measure of the utility of a reasoning path. Let’s formally define the utility score as:

$$s(z) := \pi_{\theta}(y^* | x, z) \in [0, 1]. \quad (4)$$

Then we can define the ability of the LLM as the expected utility over its reasoning distribution:

$$\mu := \mathbb{E}_{z \sim \pi_{\theta}(\cdot|x)} [s(z)] \quad (5)$$

Note that this equals to $\pi_{\theta}(y^* | x)$ because of the law of total probability. We can define the τ -good set as $\mathcal{Z}_{\tau} = \{z : s(z) \geq \tau\}$, where $\tau \geq \mu$ emphasizes above-average reasoning paths.

The current challenge is to more efficiently sample reasoning paths with higher $s(z)$. We can prove that observing y^* can increase the sampling probability of reasoning paths with higher $s(z)$. To begin with, for a specific reasoning path z , by the law of total probability and Bayes’ rule, we have

$$\begin{aligned} \mathbb{E}_{c \sim \pi_{\theta}(\cdot|x, y^*)} \pi_{\theta}(z|x, y^*, c) &= \pi_{\theta}(z | x, y^*) \\ &= \frac{\pi_{\theta}(y^* | x, z) \pi_{\theta}(z | x)}{\pi_{\theta}(y^* | x)} = \frac{s(z)}{\mu} \pi_{\theta}(z | x), \end{aligned} \quad (6)$$

where the first expression describes that the LLM thinks through the problem and reference answer jointly before generating the reasoning path. It shows that *observing y^* induces a size-biased reweighting of $\pi_{\theta}(z | x)$ by $s(z)$: high- $s(z)$ paths gain probability mass, low- $s(z)$ paths lose it.*

For any subset of reasoning paths, denoted by Z , the probability of sampling one reasoning path belonging to Z can be formulated as follows

$$\Pr(Z | x, y^*) = \Pr(Z | x) \cdot \frac{\mathbb{E}[s(z) | z \in Z, x]}{\mu} \quad (7)$$

See Appendix A.3 for detailed derivation. Consequently, we have

$$\begin{aligned} \Pr(Z | x, y^*) &\geq \Pr(Z | x) \\ \iff \mathbb{E}[s(z) | z \in Z, x] &\geq \mu. \end{aligned} \quad (8)$$

For τ -good set $\mathcal{Z}_{\tau} = \{z : s(z) \geq \tau\}$, we have

$$\frac{\Pr(\mathcal{Z}_{\tau} | x, y^*)}{\Pr(\mathcal{Z}_{\tau} | x)} \geq \frac{\tau}{\mu}, \text{ with increase if } \tau > \mu. \quad (9)$$

Therefore, conditioning on reference answer y^* amplifies above-average reasoning paths. Beyond sets, the posterior raises the expected utility over the distribution, because given equation 6 and $\mathbb{E}[s(z)^2] = \mathbb{E}[s(z)]^2 + \text{Var}(s(z))$, we can derive

$$\begin{aligned} \mathbb{E}_{z \sim \pi_{\theta}(\cdot|x, y^*)} [s(z)] &= \frac{1}{\mu} \mathbb{E}_{z \sim \pi_{\theta}(\cdot|x)} [s(z)^2] \\ &= \mu + \frac{\text{Var}(s(z))}{\mu} \geq \mu. \end{aligned} \quad (10)$$

In summary, approximate $\pi_{\theta}(z | x, y^*)$ acts as a principled target for exploring high-utility reasoning paths and thus it can be properly used to help the learning of $\pi_{\theta}(z | x)$. See Appendix A.4 for further rationality analysis.

2.3 Reference-Answer-Conditioned Variational Optimization Objective

Building on the above insight, we propose a novel variational objective, introducing an amortized posterior $\pi_{\theta}(z | x, y^*)$ to aid the learning of the prior $\pi_{\theta}(z | x)$. We start from log-transformation of the raw objective, which maximizes the utility score over the reasoning distribution:

$$\log \mathcal{J}(\theta) = \log \mathbb{E}_{z \sim \pi_{\theta}(z|x)} [\pi_{\theta}(y^* | x, z)]. \quad (11)$$

By introducing the amortized posterior $\pi_{\theta}(z | x, y^*)$ and applying Jensen’s inequality, we can derive the Evidence Lower Bound of the objective:

$$\begin{aligned} \log \mathcal{J}(\theta) &\geq \mathbb{E}_{z \sim \pi_{\theta}(z|x, y^*)} [\log \pi_{\theta}(y^* | x, z)] \\ &\quad - \mathbb{D}_{\text{KL}} [\pi_{\theta}(z|x, y^*) || \pi_{\theta}(z|x)]. \end{aligned} \quad (12)$$

See Appendix A.5 for derivation. Intuitively, the first term encourage reasoning paths that make y^* more likely. The KL term pulls the prior toward the posterior and regularizes the posterior to prevent it from collapsing onto out-of-distribution reasoning path for the prior.

To further stabilize training, we innovatively introduce a utility baseline—the expected utility under the question-only prior—to estimate reasoning utility under the answer-conditioned posterior. This provides a more informative reward signal by measuring how much the reasoning under the posterior improves over that under the question-only prior.

$$\begin{aligned} R_{\text{impr}}(z) &= \max(0, \log \pi_{\theta}(y^* | x, z) \\ &\quad - \mathbb{E}_{z' \sim \pi_{\theta}(\cdot|x)} \log \pi_{\theta}(y^* | x, z')). \end{aligned} \quad (13)$$

We clip the minimum reward to zero to stabilize training. Moreover, we replace the likelihood with length-normalized sequence likelihood. This avoids the issue where longer reference answers receive smaller rewards, which would make reference answer of different lengths incomparable. Additionally, we address a subtle but significant bias in how existing probability-reward-based methods (Yu et al., 2025b; Zhou et al., 2025b) estimate answer likelihood. These methods compute $\pi_\theta(y^* | x, z)$ by directly appending the reference answer after the reasoning path, which fails to control for the model’s internal state transition to answer generation. Our solution is to standardize this transition by inserting a cue phrase, estimating the probability as $\pi_\theta(y^* | x, z, \text{“The answer is ”})$ within the structured sequence *“The answer is y^* ”*. This prefix aligns with common linguistic patterns and establishes a consistent state for probability estimation across all reasoning paths, thereby guaranteeing the validity of the reward normalization.

Moreover, once the model has learned high-quality reasoning for a given problem, it no longer needs external guidance and should learn only from paths with higher utility than its own. Accordingly, we apply reward-based weighting to samples when estimating the KL.

$$\tilde{\mathbb{D}}_{\text{KL}} = R_{\text{impr}}(z) \cdot \mathbb{D}_{\text{KL}}[\pi_\theta(z|x, y^*) || \pi_\theta(z|x)] \quad (14)$$

Therefore, the final variational objective is:

$$\mathcal{J}_{\text{RAVR}}(\theta) = \mathbb{E}_{z \sim \pi_\theta(z|x, y^*)} [R_{\text{impr}}(z)] - \tilde{\mathbb{D}}_{\text{KL}}[\pi_\theta(z|x, y^*) || \pi_\theta(z|x)]. \quad (15)$$

We utilize GRPO to optimize the first term and jointly optimize this objective and Equation 11. To induce the two distributions $\pi(\cdot|x)$ and $\pi(\cdot|x, y^*)$ within a single LLM, we adjust the user prompt as follows.

Prompt Template for Question-Only

```
<lim_start>system
A conversation between user and assistant. The user asks a question, and the assistant solves it. The assistant first thinks about the reasoning process in the mind and then provides the user with the answer. The reasoning process is enclosed within <think></think>tags, i.e., <think>This is my reasoning. </think>This is my answer.<lim_end>
<lim_start>user
{{ question }}<lim_end>
<lim_start>assistant
<think>
```

Prompt Template for Conditioning on both Question and Reference Answer

```
# same system prompt.
<lim_start>user
Given the following question and its reference answer, your task is to produce a step-by-step explanation that logically leads to the reference answer, written in the style of a first-person think-aloud monologue. You are encouraged to draw on the reference answer for internal guidance to help structure and support your reasoning, but the final monologue must read as a genuine, first-encounter, real-time discovery, without mentioning or implying any prior access to the reference answer.
Question: {{question}}
Reference Answer: {{ground_truth}}
OUTPUT REQUIREMENTS:
1. Output ONLY the first-person, think-aloud monologue. Do not include any preface, summary, or restatement of these instructions.
2. Maintain the tone of a focused individual thinking to themselves. Avoid meta-commentary like “for the first time,” and any phrasing that reveals simulation.
3. Do not mention, imply, or hint at prior access to the Reference Answer in the monologue. Avoid phrases like “according to the answer...” or “to get to that answer...”, and any euphemism that signals foreknowledge.
4. Do not merely restate the final answer in the monologue; articulate the reasoning pathway with sufficient intermediate steps, rationale, decision points, verification, and any necessary error-correction or backtracking.
<lim_end>
<lim_start>assistant
<think>
```

3 Experiments

3.1 Performance of RA VR

We conducted experiments with Qwen3-1.7B (Qwen, 2025). We train the LLM on two datasets: CrossThink-QA (Akter et al., 2026), a general-purpose reasoning benchmark covering STEM, Economics, Social Sciences, and more, and DeepMath-103K (He et al., 2025; Liu et al., 2025), a math dataset focused on challenging problems. Our baselines span two axes: reward type (verifiable reward vs. reference-answer-probability reward) and the use of implicit curriculum learning for stabilizing training on challenging samples. We include GRPO (Shao et al., 2024), DAPO (Yu et al., 2025a), VeriFree (Zhou et al., 2025b), and RLPR (Yu et al., 2025b). GRPO and DAPO use verifiable rewards, with DAPO introducing dynamic sampling for curriculum learning; VeriFree and RLPR adopt reference-answer-probability rewards, and RLPR further applies reward standard deviation filtering to induce an adaptive curriculum. See Appendix A.6 for implementation details.

Table 1: Experiments on Qwen3-1.7B. General task represents average of GPQA-D and MMLU-Pro, and Math reasoning represents average of AIME24, AIME25, AMC23 and Minerva. The metric is the average@k accuracy, with the subscript indicating the corresponding average token length.

Training Set	Model	GPQA-D	MMLU-Pro	AIME 24	AIME 25	AMC 23	Minerva	General	Math	Average
		Avg@4	Avg@4	Avg@16	Avg@16	Avg@16	Avg@16	Task	Reasoning	
	Qwen3-1.7B	21.46 ₍₆₉₄₈₎	53.48 ₍₃₄₉₉₎	20.00 ₍₇₈₀₅₎	23.30 ₍₇₆₅₄₎	55.00 ₍₅₈₉₆₎	50.00 ₍₅₂₃₃₎	37.47	37.08	37.21
		In-domain			Out-of-domain					
CrossThink-QA	+ GRPO	34.97 ₍₄₅₈₁₎	55.18 ₍₂₅₈₇₎	25.63 ₍₇₂₇₃₎	23.54 ₍₇₂₈₃₎	64.53 ₍₄₇₉₉₎	56.92 ₍₃₂₃₁₎	45.08	42.65	43.46
	+ DAPO	35.35 ₍₅₉₉₁₎	54.75 ₍₅₁₂₉₎	23.33 ₍₇₅₅₆₎	21.88 ₍₇₄₉₅₎	62.50 ₍₅₅₀₂₎	56.80 ₍₄₁₂₉₎	45.05	41.13	42.44
	+ VeriFree	28.60 ₍₆₀₉₂₎	53.32 ₍₂₆₇₂₎	23.54 ₍₇₅₃₂₎	20.83 ₍₇₃₈₆₎	62.19 ₍₅₀₄₄₎	55.79 ₍₃₅₂₉₎	40.96	40.59	40.71
	+ RLPR	30.93 ₍₅₃₅₇₎	53.82 ₍₂₅₄₁₎	24.83 ₍₇₂₇₅₎	22.92 ₍₇₁₈₅₎	62.81 ₍₄₉₃₃₎	56.34 ₍₃₁₃₁₎	42.38	41.73	41.94
	+ RAVR (Ours)	40.91 ₍₄₁₇₇₎	55.88 ₍₂₄₃₀₎	27.92 ₍₇₂₃₆₎	23.75 ₍₇₂₈₁₎	62.03 ₍₅₀₇₃₎	58.00 ₍₃₃₂₀₎	48.39	42.92	44.75
		Out-of-domain			In-domain					
DeepMath	+ GRPO	33.21 ₍₃₉₃₂₎	54.55 ₍₁₉₃₄₎	26.67 ₍₆₉₃₅₎	21.46 ₍₆₉₀₆₎	67.03 ₍₄₆₀₂₎	56.99 ₍₂₇₆₆₎	43.88	43.04	43.32
	+ DAPO	34.84 ₍₅₂₄₂₎	55.25 ₍₂₆₈₄₎	26.67 ₍₇₂₆₃₎	23.33 ₍₇₀₇₀₎	67.50 ₍₅₀₄₃₎	57.97 ₍₃₆₃₃₎	45.05	43.87	44.26
	+ VeriFree	30.30 ₍₁₂₅₆₎	52.90 ₍₇₆₂₎	26.88 ₍₂₉₀₁₎	21.88 ₍₂₄₄₃₎	62.34 ₍₁₄₉₀₎	55.01 ₍₇₁₃₎	41.60	41.53	41.55
	+ RLPR	31.31 ₍₅₄₆₄₎	54.27 ₍₂₅₅₈₎	26.04 ₍₇₃₁₈₎	22.50 ₍₇₂₂₄₎	62.19 ₍₄₉₁₉₎	55.51 ₍₃₁₆₈₎	42.79	41.56	41.97
	+ RAVR (Ours)	34.60 ₍₄₅₆₆₎	55.50 ₍₂₅₃₈₎	29.17 ₍₆₅₉₃₎	22.71 ₍₆₈₃₄₎	69.69 ₍₄₆₁₅₎	58.43 ₍₃₂₆₁₎	45.05	45.00	45.02

The results in Table 1 demonstrate the superior performance and generalizability of **RAVR**. When trained on CrossThink-QA, **RAVR** achieves the highest average score, leading on in-domain general tasks and exhibiting the best out-of-domain performance on math benchmarks. When trained on DeepMath, **RAVR** again attains the highest overall average and shows strong generalizability to out-of-domain general tasks. Overall, **RAVR** consistently delivers state-of-the-art results across training settings and evaluation benchmarks. We further analyze the relationship between reasoning length and performance. All RL-based methods substantially outperform the vanilla model while typically using fewer tokens, and **RAVR** is generally more

token-efficient: when trained on CrossThink-QA, it uses fewer tokens than GRPO on 4 of 6 datasets; when trained on DeepMath, it matches GRPO’s token budget and is more efficient than DAPO. This suggests a favorable length regime in which reasoning is sufficiently detailed but not verbose. **RAVR** tends to remain in this regime, likely because answer-conditioned reasoning prunes unnecessary steps and encourages compact yet effective reasoning during training.

3.2 Generalizability Analysis of RAVR

To further assess the generalizability of **RAVR**, we compare it with GRPO across a broader set of datasets and models. As shown in Figure 2, the

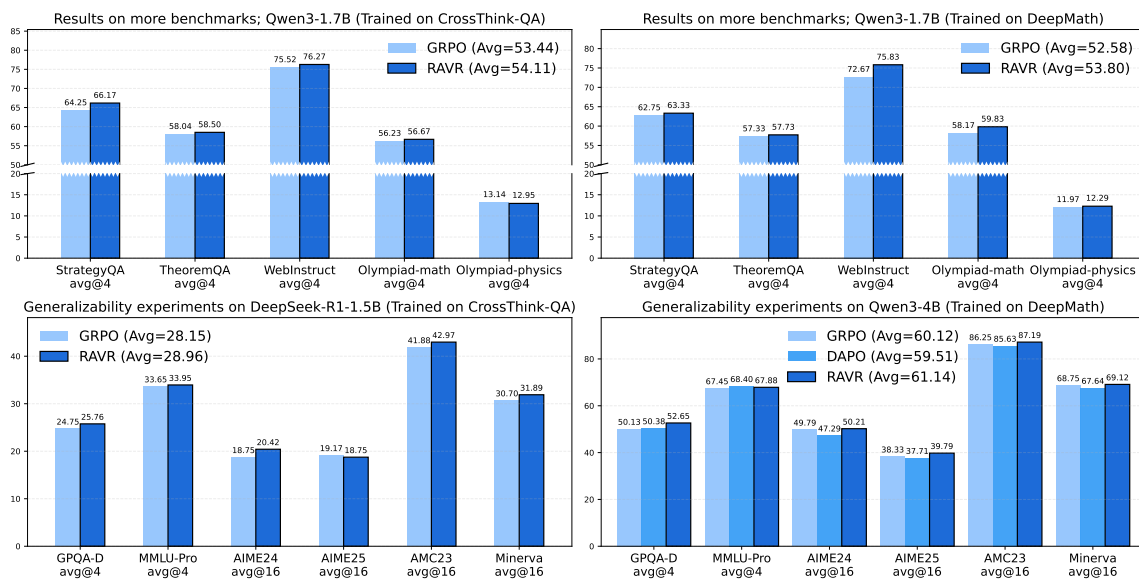
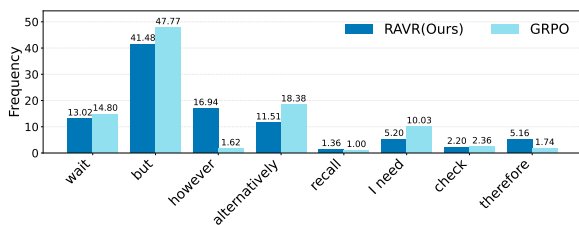


Figure 2: Generalizability analysis.

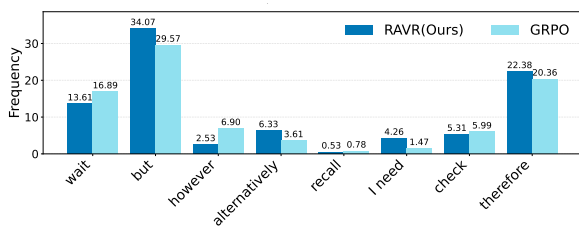
advantages of **RAVR** consistently persist across settings. We first consider benchmarks requiring stronger compositional reasoning beyond multiple-choice and standard math problems, like TheoremQA, which evaluates formal theorem application on challenging science problems. We then evaluate DeepSeek-R1-1.5B trained on CrossThink-QA. **RAVR** outperforms GRPO on 5 of 6 datasets, improving the overall average accuracy by about 2.9%. We also study a larger model, Qwen3-4B. As it already achieves nearly 90% accuracy on CrossThink-QA, we instead train on the more challenging DeepMath dataset and include DAPO as an additional strong baseline. **RAVR** consistently surpasses GRPO and outperforms DAPO on 5 of 6 datasets, with average gains of about 1.7% and 2.7%, respectively. See Appendix A.7 for additional generalizability analysis of Qwen3-8B and a more challenging reasoning scenario.

3.3 Reasoning Behavior of RAVR

Figure 3 shows the gap in thinking style of the model trained with **RAVR** and GRPO via frequencies of discourse markers as proxies for cognitive moves (Wang et al., 2025; Bogdan et al., 2025).



(a) Trained on CrossThink-QA; Tested on GPQA-Diamond.



(b) Trained on DeepMath; Tested on AMC23.

Figure 3: Comparison of reasoning behaviors during thinking. Words on x-axis are those frequently used during thinking. Y-axis represents their average frequency per response. See Appendix A.8 for more results.

The key reasoning features of **RAVR** are observed as follows: (1) *Fewer wait*. This indicates reduced dithering, which is consistent with stronger problem-solving competence and with answer-conditioned paths that tend to terminate once the

reference answer is reached and avoid overthinking. (2) *More therefore*. This suggests firmer result consolidation: the model more actively reviews preceding steps and commits to a conclusion. (3) *More recall in knowledge QA*. This necessitates extensive retrieval to ensure sufficient recall for Knowledge QA, rather than jumping straight to an answer. (4) *More alternatively and I need in math*. These markers reflect greater divergent exploration and explicit planning before committing to a solution path—desirable for multi-step problem solving. (5) *Task-adaptive contrast*. In English discourse, *but* and *however* both serve as contrastive markers, yet they differ in usage. *But* is more casual and often signals a local correction or small-scale turn within a sentence, whereas *however* is more formal and typically marks a global contrast or structured shift across sentences. Multiple-choice task exhibits more *however* (global comparison across options), while math shows more *but* (incremental corrections within derivations), suggesting an adjustment of *contrast granularity* to task demands rather than a fixed stylistic habit. The reasoning behavior shifts indicate a more interpretable, problem-adaptive reasoning process that **RAVR** achieves.

3.4 Learning Dynamics of RAVR

To further investigate how introducing answers enhances sampling efficiency, we compare our model against GRPO with larger rollout group sizes.

Figure 4 demonstrates the superior sampling efficiency of **RAVR**. While the performance of GRPO scales with a larger group size, **RAVR** achieves a better or comparable performance with a smaller group size of 8. This finding provides strong evidence that **RAVR** markedly enhances the sampling efficiency of high-quality reasoning paths. Furthermore, the smoother learning curve for **RAVR** indicates improved learning stability.

Moreover, Figure 5 shows that the KL divergence between $\pi_\theta(z|x, y^*)$ and $\pi_\theta(z|x)$ first fluctuates but then gradually decreases. This indicates that the capability of producing high-utility reasoning paths under the posterior is transferred to the question-only setting as expected, and that the language style of the generated reasoning increasingly aligns with that of reasoning when question-only. Additionally, it’s observed that while the prior reasoning utility becomes better, the posterior keeps a stable utility gain, which ensures the continual learning. See Appendix A.9 for more dynamics.

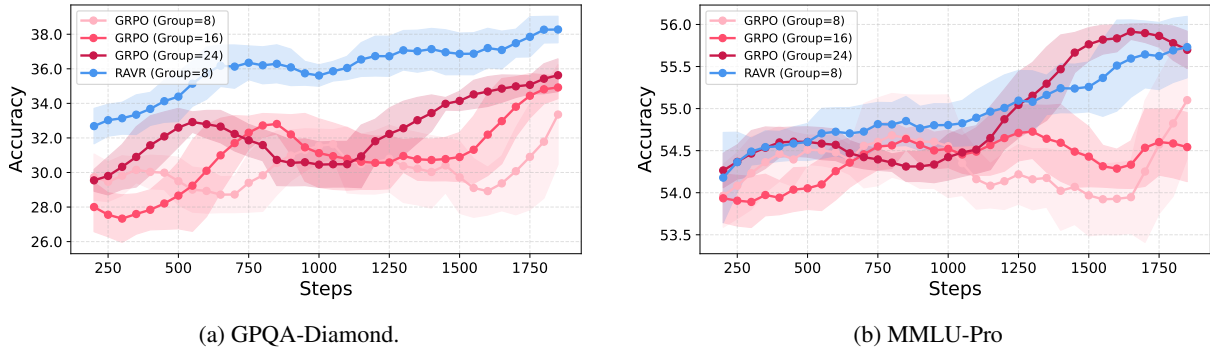


Figure 4: Comparison with GRPO across different rollout group sizes. It indicates that **RAVR** markedly enhances the sampling efficiency of high-quality reasoning, thereby improving learning stability and efficiency.

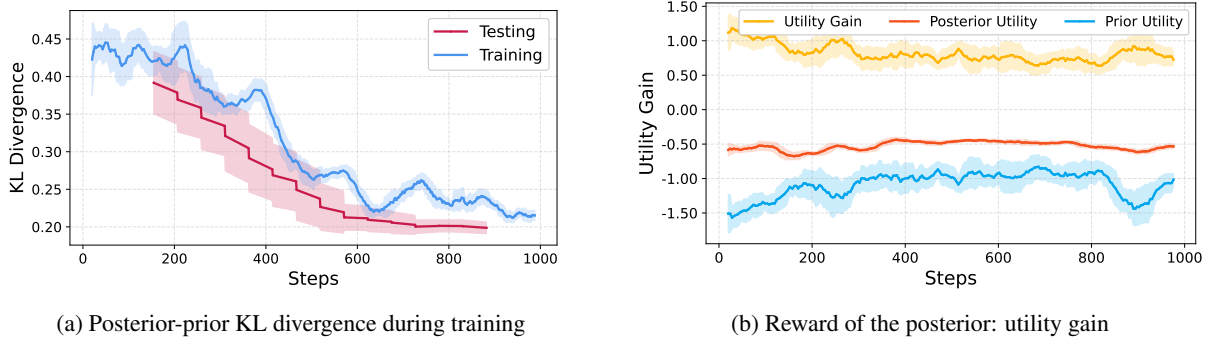


Figure 5: KL Dynamics. Training on CrossThink-QA, testing on GPQA and MMLU-pro.

3.5 Ablation Study of RAVR

To validate the design choices of **RAVR**, we conduct ablation studies on its key components, demonstrating that each contributes indispensably to the overall stability and performance. (1) *Main variational reasoning objective*. Maximizing the utility of answer-conditioned reasoning is critical: removing this component makes **RAVR** consistently underperform GRPO across both datasets. Without it, the KL term can pull the question-only prior toward a diffused posterior, thereby hindering effective prior learning. Removing the KL term,

in turn, reduces training to an auxiliary reasoning task that produces reasoning given both the question and the answer. On the math dataset, we even observe slight gains in this setting, likely because the answer adds little beyond flagging errors and avoiding overthinking, making posterior and prior reasoning largely similar. (2) *Utility baseline and posterior instruction*. Both strategies are important. Without the prior-based utility baseline, posterior reasoning is no longer required to surpass prior reasoning, and training tends to drift toward “least bad” patterns rather than genuinely informa-

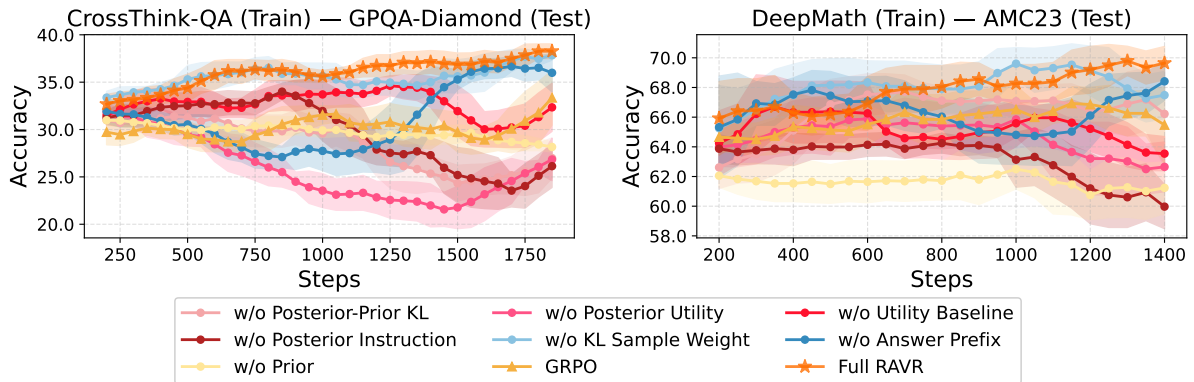


Figure 6: Ablation study.

tive ones, resulting in only marginal improvement. Omitting the instruction to produce first-person, think-aloud monologues also destabilizes training, likely because the resulting language-style shift disrupts the model’s original reasoning patterns. (3) *KL sample weighting and answer prefix*. These stability-oriented strategies play distinct roles. KL reweighting has little effect on peak accuracy, but removing it leads to higher variance during training. The answer prefix—absent from prior probability-reward methods—helps stabilize early training and enables faster, more consistent convergence. (4) *Prior objective*. Retaining the question-only reasoning objective is vital, especially for math, since mismatches between training (posterior) and inference (prior) can degrade the performance of a probabilistic model. Overall, these ablations validate the full design of **RAVR**, showing that it achieves better stability and performance.

4 Related Work

Earlier works have shown that revealing the correct answer to the LLM can help in generating useful reasoning, which can then be used for supervised fine-tuning (SFT). For instance, [Zelikman et al. \(2022\)](#) introduced STaR, a bootstrapping loop that first lets the model attempt each problem on its own, and if it fails, then provides the ground-truth and asks the model to solve again. Recently, [Zhou et al. \(2025a\)](#) expands STaR into the RL setting, which uses STaR to generate positive samples for preference optimization, such as DPO ([Rafailov et al., 2023](#)). However, the resulting reasoning is usually filtered primarily by whether it ends with correct answer. This leaves limited control over the quality of the intermediate reasoning itself, such as whether it reflects genuine problem solving rather than post-hoc rationalization or answer leakage. Recently, [Cetin et al. \(2025\)](#) proposes an approach to explicitly train the ability of the LLM to generate reasoning with reference answer, thereby yielding better reasoning paths for SFT. However, this work requires complete step-by-step answer of a problem and cannot handle scenarios where detailed answers are unavailable.

Moreover, while these methods demonstrate the value of answer-guided exploration, they typically rely on multi-stage pipelines that alternate between answer-guided exploration and SFT, rather than optimizing answer-guided exploration and question-only reasoning under a unified, end-to-end objec-

tive within a single LLM. Additionally, the reference answer is often treated merely as a prompt for immediate reasoning generation, overlooking the LLM’s capacity to reason using the reference answer for explanation derivation. In contrast, **RAVR** explicitly unifies answer-conditioned exploration with question-only reasoning within a single LLM in an end-to-end manner. It enables the LLM to leverage reference answers to reason about *why this is the answer*, rather than treating them merely as surface-level hints. Furthermore, **RAVR** fully exploits the LLM’s reasoning capabilities, requiring no complete reference solutions—a single word or phrase suffices.

5 Conclusion

RAVR is an end-to-end framework that makes the LLM think *why is this the answer* to help its learning of *what is the answer* via a novel variational reinforcement learning objective. This is the first work that leverages the reasoning ability of LLMs to use the reference answer to derive high-quality reasoning. We further innovatively propose simple but effective strategies, namely the reasoning reward baseline, reward-based sample weight in KL calculation and the answer prefix trick, to enhance the training stability and performance.

Limitations

Our study primarily focuses on reasoning tasks with relatively well-defined reference answers, which is common in RL-based LLM reasoning. This evaluation setting does not fully cover more open-ended scenarios, including long-form writing or free-form explanation, where reference answers typically provide richer and softer supervision and where it is generally more difficult for LLMs to obtain high-quality reasoning solely from the question and their inherent capabilities. Importantly, this limitation reflects the scope of our current evaluation rather than a fundamental restriction of the proposed framework, and extending the answer-conditioned variational objective to more open-ended forms of reference supervision is a promising direction for future work. In addition, due to computational resource constraints, our experiments are limited to small- and medium-scale models. While the observed improvements are consistent across multiple benchmarks and model architectures, we leave a systematic investigation of larger model sizes to future work.

References

- Syeda Nahida Akter, Shrimai Prabhumoye, Matvei Novikov, Seungju Han, Ying Lin, Evelina Bakhturina, Eric Nyberg, Yejin Choi, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. 2026. [Nemotron-crossthink: Scaling self-learning beyond math reasoning](#). In *Proceedings of the 19th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2026 - Volume 1: Long Papers, Rabat, Morocco, March 24-29, 2026*, pages 984–1002. Association for Computational Linguistics.
- Paul C Bogdan, Uzay Macar, Neel Nanda, and Arthur Conmy. 2025. Thought anchors: Which llm reasoning steps matter? *arXiv preprint arXiv:2506.19143*.
- Edoardo Cetin, Tianyu Zhao, and Yujin Tang. 2025. Reinforcement learning teachers of test time scaling. *arXiv preprint arXiv:2506.08388*.
- Micheline TH Chi, Miriam Bassok, Matthew W Lewis, Peter Reimann, and Robert Glaser. 1989. Self-explanations: How students study and use examples in learning to solve problems. *Cognitive science*, 13(2):145–182.
- DeepSeek-AI. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *CoRR*, abs/2501.12948.
- Samuel Gershman and Noah Goodman. 2014. Amortized inference in probabilistic reasoning. In *Proceedings of the annual meeting of the cognitive science society*, volume 36.
- Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu Chen, Yue Wang, Linfeng Song, Dian Yu, Zhenwen Liang, Wenxuan Wang, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. 2025. [Deepmath-103k: A large-scale, challenging, decontaminated, and verifiable mathematical dataset for advancing reasoning](#). *CoRR*, abs/2504.11456.
- Shulin Huang, Linyi Yang, Yan Song, Shuang Chen, Leyang Cui, Ziyu Wan, Qingcheng Zeng, Ying Wen, Kun Shao, Weinan Zhang, Jun Wang, and Yue Zhang. 2025. [Thinkbench: Dynamic out-of-distribution evaluation for robust LLM reasoning](#). *CoRR*, abs/2502.16268.
- Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- Zongxia Li, Yapei Chang, Yuhang Zhou, Xiyang Wu, Zichao Liang, Yoo Yeon Sung, and Jordan Lee Boyd-Graber. 2025. Semantically-aware rewards for open-ended r1 training in free-form generation. *arXiv preprint arXiv:2506.15068*.
- Zihe Liu, Jiashun Liu, Yancheng He, Weixun Wang, Jiaheng Liu, Ling Pan, Xinyu Hu, Shaopan Xiong, Ju Huang, Jian Hu, Shengyi Huang, Siran Yang, Jiamang Wang, Wenbo Su, and Bo Zheng. 2025. [Part I: tricks or traps? A deep dive into RL for LLM reasoning](#). *CoRR*, abs/2508.08221.
- OpenAI. 2024. [Openai o1 system card](#). *CoRR*, abs/2412.16720.
- OpenAI. 2025. [gpt-oss-120b & gpt-oss-20b model card](#). *CoRR*, abs/2508.10925.
- Qwen. 2025. [Qwen3 technical report](#). *CoRR*, abs/2505.09388.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741.
- John Schulman. 2020. [Approximating kl divergence](#).
- Murray Shanahan, Kyle McDonell, and Laria Reynolds. 2023. Role play with large language models. *Nature*, 623(7987):493–498.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. [Deepseekmath: Pushing the limits of mathematical reasoning in open language models](#). *CoRR*, abs/2402.03300.
- Chenlong Wang, Yuaning Feng, Dongping Chen, Zhaoyang Chu, Ranjay Krishna, and Tianyi Zhou. 2025. Wait, we don't need to "wait"! removing thinking tokens improves reasoning efficiency. *arXiv preprint arXiv:2506.08343*.
- Yifei Xu, Tusher Chakraborty, Srinagesh Sharma, Leonardo Nunes, Emre Kiciman, Songwu Lu, and Ranveer Chandra. 2025. Direct reasoning optimization: LLMs can reward and refine their own reasoning for open-ended tasks. *arXiv preprint arXiv:2506.13351*.
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, and 16 others. 2025a. [DAPO: an open-source LLM reinforcement learning system at scale](#). *CoRR*, abs/2503.14476.
- Tianyu Yu, Bo Ji, Shouli Wang, Shu Yao, Zefan Wang, Ganqu Cui, Lifan Yuan, Ning Ding, Yuan Yao, Zhiyuan Liu, Maosong Sun, and Tat-Seng Chua. 2025b. [RLPR: extrapolating RLVR to general domains without verifiers](#). *CoRR*, abs/2506.18254.
- Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Shiji Song, and Gao Huang. 2025. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model? *arXiv preprint arXiv:2504.13837*.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. 2022. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488.

Kaiyan Zhang, Yuxin Zuo, Bingxiang He, Youbang Sun, Runze Liu, Che Jiang, Yuchen Fan, Kai Tian, Guoli Jia, Pengfei Li, Yu Fu, Xingtai Lv, Yuchen Zhang, Sihang Zeng, Shang Qu, Haozhan Li, Shijie Wang, Yuru Wang, Xinwei Long, and 20 others. 2025a. [A survey of reinforcement learning for large reasoning models](#). *CoRR*, abs/2509.08827.

Kaiyi Zhang, Ang Lv, Jinpeng Li, Yongbo Wang, Feng Wang, Haoyuan Hu, and Rui Yan. 2025b. [Stephint: Multi-level stepwise hints enhance reinforcement learning to reason](#). *arXiv preprint arXiv:2507.02841*.

Ruiyang Zhou, Shuoze Li, Amy Zhang, and Liu Leqi. 2025a. [Expo: Unlocking hard reasoning with self-explanation-guided reinforcement learning](#). *arXiv preprint arXiv:2507.02834*.

Xiangxin Zhou, Zichen Liu, Anya Sims, Haonan Wang, Tianyu Pang, Chongxuan Li, Liang Wang, Min Lin, and Chao Du. 2025b. [Reinforcing general reasoning without verifiers](#). *arXiv preprint arXiv:2505.21493*.

A Appendix

A.1 Information About Use of AI Assistants

In the course of this work, AI models were employed as auxiliary tools in the following ways:

1. **Figure generation** — The model *Gemini-2.5-Pro* (with *Nano Banana*) was used to assist in the creation of Figure 1.
2. **Text editing** — *GPT-5* was utilized for grammar checking, paraphrasing, and shortening of sections of the manuscript.
3. **Code assistance** — *GPT-5* was used to aid in code implementation, particularly for visualization-related scripts.

A.2 Motivation Experiment

We used Qwen3-1.7B with thinking mode enabled and ran inference using the officially recommended settings: Temperature=0.6, TopP=0.95, TopK=20; the maximum completion tokens were set to 8192.

One experimental case is as follows. The question asked: *Which renewable energy source has the highest capacity factor, making it the most reliable for consistent power generation?* The correct answer is *geothermal energy*. When only the question was provided, the model predicted *hydroelectric energy*, overestimating hydro’s capacity factor by assuming that water is always flowing and overlooking its dependence on water availability—some plants are used for peak power, so actual average can be lower. It also failed to analyze geothermal

in depth after deeming *hydropower* the correct answer and arrived at an incorrect, intuition-driven estimate of geothermal’s capacity factor.

When the reference answer was provided, the model noted that *hydroelectric* output can vary with location and season, and it carefully reasoned that geothermal plants can operate continuously as long as there is a viable geothermal resource; once the plant is running, it can produce steady power. It also gave more accurate estimates of the capacity factors of hydroelectric and geothermal sources. As a result, GPT-5 evaluated the derived reasoning as: *It correctly defines capacity factor, compares sources, and concludes geothermal has the highest capacity factor and reliability; minor oversimplifications (e.g., biomass) do not affect the correct conclusion. The reasoning is self-contained, does not reference the provided answer, and shows no signs of relying on it.*

Overall, *when given the reference answer, the model did not collapse into answer-conforming shortcuts. On the contrary, the presence of the answer promoted a more detailed discussion of both the initially uncertain (incorrect) option and the correct one.*

Case Study of Motivation Experiment

A. Question; Reference Answer of This Question is D.

Which of the following renewable energy sources has the highest capacity factor, making it the most reliable for consistent power generation?

- A: Solar photovoltaic (PV) energy
- B: Wind energy
- C: Hydroelectric energy
- D: Geothermal energy
- E: Biomass energy
- F: Tidal energy
- G: Wave energy
- H: Solar thermal energy
- I: Biofuels
- J: Hydrogen fuel cells

Please think step by step and output the final answer within `\boxed{}`.

B. Reasoning and Final Response When Question-only; Using Prompt in Section 2.3.

Okay, let’s try to figure out which renewable energy source has the highest capacity factor. First, I need to remember what capacity factor is. From what I recall, capacity factor is the ratio of actual output over the maximum possible output. So, it’s like how much of the potential energy a power plant can produce in a given time. A higher capacity factor means the energy source is more reliable because it produces more consistently.

Now, looking at the options: A to J. Let me go through each one. Starting with Solar photovoltaic (PV) energy. I think solar PV depends on sunlight, which is variable. Even on cloudy days, it might not produce as much. So maybe its capacity factor is lower compared to others. But I’m not sure. I remember that solar panels have a capacity factor around 15-20% for PV, depending on location and weather. Maybe that’s not the highest.

Wind energy. Wind power is also variable because wind doesn’t blow all the time. The capacity factor for wind is around 30-40% in some places, but it can vary. For example, in areas with consistent wind, like coastal regions, it might be higher. But I think wind’s capacity factor is lower than hydroelectric. Wait, but maybe in some places it’s higher?

Hydroelectric energy. Hydroelectric plants have a high capacity factor

because they use the water flow from rivers, which is more consistent. They can generate electricity 24/7 as long as there's water flow. I think hydroelectric has a capacity factor around 80-90% in many places. That seems really high. But wait, does that apply to all hydroelectric plants? Maybe some are more efficient, but generally, hydro is pretty reliable.

Geothermal energy. Geothermal plants can be pretty reliable because they use the Earth's internal heat. But their capacity factor depends on the location. For example, in places with consistent geothermal activity, like Iceland or the US, they might have a high capacity factor. But I think it's around 50-60% maybe? Not sure. Maybe lower than hydro.

Biomass energy. Biomass is from organic materials. It's variable because it depends on the availability of biomass. For example, if there's a lot of crops or waste, it can be more reliable. But the capacity factor for biomass is usually lower than hydro. Maybe around 20-30%?

Tidal energy. Tidal energy depends on the tides, which are predictable. So maybe it has a higher capacity factor? But I think it's lower than hydro because tides are less consistent than river flows. But I'm not sure. Tidal energy's capacity factor might be around 50-60%, maybe? Wave energy. Similar to tidal, but wave energy is more variable. The capacity factor for wave energy is probably lower than tidal. Maybe around 20-30%?

Solar thermal energy. Solar thermal uses sunlight to generate heat, which then produces electricity. But it's similar to solar PV in that it depends on sunlight. So maybe similar capacity factor as PV, around 15-20%.

Biofuels. Biofuels are made from organic matter, but they can be less reliable because they depend on agricultural inputs. Their capacity factor might be lower than hydro. Maybe around 20-30%?

Hydrogen fuel cells. Hydrogen is produced via electrolysis, which is usually from renewable sources. But hydrogen fuel cells themselves are not a power source. Wait, the question is about renewable energy sources. So hydrogen fuel cells are not a source of energy but a method of storing energy. So maybe not part of the options here. Wait, the options include "Hydrogen fuel cells" as an option. But the question is about renewable energy sources, so maybe hydrogen is considered a renewable source, but the fuel cell is the storage. So maybe the answer is not here. But the question says "renewable energy sources", so hydrogen fuel cells are not a source, but a method. So maybe not part of the options. Wait, but the options are all renewable energy sources. Wait, the options are: A to J. So maybe hydrogen fuel cells are considered a renewable energy source? Or is it not? Hmm. Maybe not. So perhaps the answer is hydroelectric.

So, going back. Hydroelectric has the highest capacity factor, like 80-90% in some places. Wind is around 30-40%, solar PV 15-20%, geothermal 50-60%, tidal 50-60%, wave 20-30%, biomass 20-30%, biofuels 20-30%, hydrogen fuel cells... not sure. So the highest is hydroelectric. So the answer would be C: Hydroelectric energy.

But wait, I need to check if there's any other sources with higher capacity factors. For example, in some regions, wind might have higher capacity factors. But globally, hydroelectric is the highest. Also, the question says "the most reliable for consistent power generation". Hydroelectric is the most reliable because it's based on water flow, which is more consistent than wind or solar. So yes, I think the answer is C.

</think>

To determine which renewable energy source has the highest capacity factor, we must first understand what capacity factor means. It is defined as the ratio of actual energy output over the maximum possible output of a power plant, typically expressed as a percentage. A higher capacity factor indicates that a power source is more reliable and consistent in producing energy over time.

Key Renewable Energy Sources and Their Capacity Factors

Let's evaluate each option:

- A: Solar photovoltaic (PV) energy
 - Capacity factor: Typically 15–20%
 - Variability: Dependent on sunlight, which is limited by weather and time of day.
- B: Wind energy
 - Capacity factor: 30–40% in some locations, but can drop significantly in low-wind areas.
 - Variability: Wind is unpredictable, so capacity factor is lower than hydroelectric.
- C: Hydroelectric energy
 - Capacity factor: 80–90% in many regions, especially where there is consistent water flow (e.g., rivers, dams).
 - Reliability: Water flow is predictable and consistent, making it one of the most reliable sources.
- D: Geothermal energy
 - Capacity factor: 50–60% in regions with stable geothermal activity.
 - Variability: Relatively stable but not as consistent as hydroelectric.
- E: Biomass energy

- Capacity factor: 20–30%.
 - Variability: Relies on biomass availability, which can be limited.
- F: Tidal energy
 - Capacity factor: 50–60%.
 - Variability: Tides are predictable but less consistent than river flows.
 - G: Wave energy
 - Capacity factor: 20–30%.
 - Variability: Less predictable and more variable than tidal or hydroelectric.
 - H: Solar thermal energy
 - Capacity factor: 15–20%.
 - Similar to solar PV in terms of variability.
 - I: Biofuels
 - Capacity factor: 20–30%.
 - Variability: Dependent on agricultural inputs.
 - J: Hydrogen fuel cells
 - Not a source of energy but a form of energy storage.
 - Not considered a renewable energy source in this context.

Conclusion

Among all the listed options, hydroelectric energy stands out as the most reliable and consistent source of power due to its predictable and stable energy output. It has the highest capacity factor, typically 80–90%, and is not subject to the same variability as solar or wind energy.

Final Answer

\boxed{C}

C. Deriving Valid Reasoning with Reference Answer; Using Prompt in Section 2.3.

Okay, let's tackle this question. The user is asking which renewable energy source has the highest capacity factor, making it the most reliable for consistent power generation. The options are from A to J, and the reference answer is D, which is geothermal energy. But I need to figure this out step by step without referencing the answer.

First, I need to recall what capacity factor means. Capacity factor is the ratio of actual output to the maximum possible output over a given period. A higher capacity factor means the energy source can produce more consistent power. So, the question is about which source can produce the most energy on average, even if it's not always producing at maximum capacity.

Let me go through each option one by one. Starting with solar PV. I know that solar energy depends on sunlight, which is intermittent. It's not always sunny, so the capacity factor is lower. For example, solar PV might have a capacity factor around 20-30% depending on location and weather.

Next, wind energy. Wind is also intermittent. The capacity factor for wind varies, but it's usually around 30-40% in some regions. However, wind turbines can be more variable because wind speeds fluctuate. But I remember that in some areas, wind farms have higher capacity factors. But I need to check if that's the case.

Hydroelectric energy. This is more reliable. Hydropower plants can generate electricity when there's water flow, which is typically consistent. However, the capacity factor for hydro is around 40-60% depending on the location and the season. But wait, some hydro plants have higher capacity factors because they can store water and release it when needed. But I need to confirm if that's accurate.

Geothermal energy. This is a bit more unique. Geothermal plants can operate continuously as long as there's geothermal activity. The capacity factor for geothermal is usually high, maybe around 60-70%? Because once the plant is running, it can produce consistently. But I'm not sure about the exact numbers.

Biomass energy. This is a bit tricky. Biomass can be intermittent because it depends on agricultural waste or organic materials. The capacity factor might be lower compared to others. Also, biomass can be more variable depending on the feedstock and availability.

Tidal energy. Tidal energy is predictable because tides are regular. But the capacity factor might be lower than hydro because the energy is tied to the tides, which are not as consistent as the river flow. However, some tidal plants have high capacity factors.

Wave energy. Similar to tidal, but wave energy is more variable. The capacity factor might be lower than tidal, but I'm not sure.

Solar thermal energy. This is a type of solar power that uses concentrated sunlight. It's similar to solar PV but can have higher efficiency. However, the capacity factor might be similar to solar PV because it's still dependent on sunlight availability.

Biofuels. These are derived from organic matter. They can be intermittent if the feedstock is not available, but in some cases, they can be more consistent. However, the capacity factor might be lower than others.

Hydrogen fuel cells. These are a type of energy storage, not a power generation source. So they might not have a high capacity factor because they are used for storage, not generation. Wait, but the question is about renewable energy sources for power generation.

Hydrogen fuel cells are more about storage, so maybe they are not considered here. So maybe they are not in the options. Wait, the options include J: Hydrogen fuel cells. But the question is about renewable energy sources for power generation. So hydrogen fuel cells might be considered as a source of energy, but their capacity factor is lower because they require storage. So maybe they are not the best answer.

So, going back. The options are A to J. *The reference answer is D: geothermal. But I need to figure out why.*

Hydroelectric is a good candidate because it's consistent. But geothermal is even better because it's a continuous source. However, I need to check if geothermal is the highest capacity factor.

Another point: geothermal plants can operate 24/7, but they require specific geological conditions. So their capacity factor is high. For example, in some regions, geothermal plants have a capacity factor over 80-90%? Wait, that seems high. But I think that's correct. Because once the plant is running, it can produce consistently. So geothermal has a very high capacity factor.

So the reasoning would be: solar PV is intermittent, wind has variable output, hydro is consistent but maybe not as high as geothermal. Geothermal is the most reliable because it's continuous. So the answer is D.

But I need to make sure there's no other option with a higher capacity factor. Let me check again.

Tidal energy is predictable but maybe not as high as geothermal. Wave energy is similar. Solar thermal is similar to solar PV. Biofuels and hydrogen are less reliable. So geothermal is the best.

Therefore, the answer is D.

</think>

I begin by defining the term "capacity factor" as the ratio of actual energy output to the maximum possible output over a given period. This metric is crucial for assessing the reliability of a renewable energy source. My goal is to determine which option provides the highest capacity factor, ensuring consistent power generation.

Starting with solar photovoltaic (PV) energy, I recognize that its output depends on sunlight availability, which is intermittent. This leads to a lower capacity factor, typically around 20-30%, depending on geographic and meteorological conditions. I note that solar PV is not a constant source, making it less reliable for consistent power generation.

Next, wind energy. While wind is a viable option, its capacity factor is influenced by variable wind speeds. In regions with consistent wind patterns, the capacity factor can reach 30-40%, but this is still less than hydroelectric sources. I consider that wind's variability limits its reliability compared to more stable energy sources.

Hydroelectric energy stands out due to its dependence on water flow, which is influenced by seasonal and climatic factors. However, hydroelectric plants often have higher capacity factors, typically ranging from 40-60%, especially in regions with consistent river flow. I acknowledge that while hydro is reliable, its capacity factor is lower than geothermal.

Geothermal energy is unique in that it operates continuously as long as there is geothermal activity. This ensures a high capacity factor, often exceeding 60-70%, depending on the site. The consistency of geothermal energy is unmatched by other sources, making it a strong candidate for reliability.

Biomass energy and hydrogen fuel cells are less reliable due to their dependence on feedstock availability and storage. Biomass is intermittent, and hydrogen fuel cells require storage infrastructure, which complicates their capacity factor. These options are not as consistent as geothermal.

Tidal and wave energy are predictable but have lower capacity factors compared to hydroelectric. Tidal energy, for example, is limited by the periodicity of tides, which results in a lower capacity factor than hydroelectric.

I reassess the options, considering the variability of each energy source. While hydroelectric and geothermal both have high capacity factors, geothermal's continuous operation and stable output make it the most reliable. I verify that geothermal's capacity factor is indeed the highest among the listed options, ensuring that my reasoning is comprehensive and accurate.

Thus, the most reliable renewable energy source for consistent power generation is geothermal energy.

D. GPT-5 Judgement.

Prompt

Please evaluate the following reasoning process on two aspects:

Question: {{question}}

Reasoning Process: {{derived reasoning process}}

Reference Answer: {{reference answer}}

Please evaluate separately:

1. Correctness Assessment: Does this reasoning process correctly lead to the answer?

- Judge whether the reasoning logic is correct

- Judge whether each step is reasonable and coherent

- Judge whether the final conclusion is correct

2. Independence Assessment: Does this reasoning process not rely on the reference answer?

- Judge whether there are obvious traces of using the reference answer

- Judge whether there are statements like "I know the answer is X"

Please answer in the following format strictly:

Correctness Assessment: Yes/No

Correctness Explanation: one or two concise sentences

Independence Assessment: Yes/No

Independence Explanation: one or two concise sentences

Judgement

Correctness Assessment: Correctness Assessment: Yes

Correctness Explanation: It correctly defines capacity factor, compares sources, and concludes geothermal has the highest capacity factor and reliability; minor oversimplifications (e.g., biomass) do not affect the correct conclusion.

Independence Assessment: Yes

Independence Explanation: The reasoning is self-contained, does not reference the provided answer, and shows no signs of relying on it.

A.3 Decomposition and Rewriting of

$$\Pr(Z | x, y^*)$$

x denotes the input problem and y^* denotes the reference answer. \mathcal{Z} is the (countable) set of all reasoning paths; $Z \subseteq \mathcal{Z}$ is a subset of reasoning paths. $\pi_\theta(z | x)$ is the model's distribution over reasoning paths given x ; $\pi_\theta(y^* | x, z)$ is the conditional distribution over reference answer given (x, z) . Define

$$s(z) := \pi_\theta(y^* | x, z), \quad \mu := \pi_\theta(y^* | x).$$

The posterior over paths given (x, y^*) induces a probability on the event Z , i.e., sampling a reasoning path that belongs to Z :

$$\Pr(Z | x, y^*) = \sum_{z \in Z} \pi_\theta(z | x, y^*). \quad (16)$$

Apply Bayes' rule, for each $z \in \mathcal{Z}$, we have

$$\begin{aligned} \pi_\theta(z | x, y^*) &= \frac{\pi_\theta(y^* | x, z) \pi_\theta(z | x)}{\pi_\theta(y^* | x)} \\ &= \frac{s(z) \pi_\theta(z | x)}{\mu} \quad (\text{assuming } \mu > 0). \end{aligned} \quad (17)$$

Substituting equation 17 into equation 16 yields

$$\begin{aligned} \Pr(Z | x, y^*) &= \sum_{z \in Z} \frac{s(z) \pi_\theta(z | x)}{\mu} \\ &= \frac{1}{\mu} \sum_{z \in Z} s(z) \pi_\theta(z | x). \end{aligned} \quad (18)$$

For a discrete space, the conditional expectation of $s(z)$ given $z \in Z$ (under $\pi_\theta(\cdot | x)$) is

$$\mathbb{E}[s(z) | z \in Z, x] = \sum_{z \in Z} s(z) \frac{\pi_\theta(z | x)}{\Pr(Z | x)}, \quad (19)$$

defined when $\Pr(Z | x) > 0$,

where by definition,

$$\Pr(Z | x) = \sum_{z \in Z} \pi_{\theta}(z | x). \quad (20)$$

Multiplying both sides of equation 19 by $\Pr(Z | x)$ gives

$$\sum_{z \in Z} s(z) \pi_{\theta}(z | x) = \Pr(Z | x) \cdot \mathbb{E}[s(z) | z \in Z, x]. \quad (21)$$

Substituting equation 21 into equation 18 yields the desired factorization:

$$\Pr(Z | x, y^*) = \Pr(Z | x) \cdot \frac{\mathbb{E}[s(z) | z \in Z, x]}{\mu}. \quad (22)$$

A.4 Rationality Analysis of Answer-Conditioned Reasoning

We further analyze the quality of answer-conditioned reasoning. In this analysis, we randomly sample 500 examples from the CrossThink-QA dataset. For each example, we enable the Qwen3-1.7’s thinking mode and generate 8 question-only prior reasoning paths. We then retain only those questions for which none of the 8 prior reasoning paths produces the correct answer, yielding a set of 85 challenging questions in total. For each of these 85 questions, the model also produces up to 8 answer-conditioned posterior reasoning paths.

The answer-conditioned reasoning paths can be categorized into four types based on whether the

reasoning is logically valid (i.e., free of hallucinations or reward hacking) and whether it is independent of the reference answer (i.e., does not indicate prior access to the reference answer): **valid and independent**, **valid and dependent**, **invalid and independent**, and **invalid and dependent**. We use GPT-5.1 to annotate each reasoning path with the same prompt in the above case study. Valid and independent reasoning paths appear on 62.4% of examples and account for 32.1% of answer-conditioned trajectories per example on average. For each reasoning path we measure the log-probability of the ground-truth answer, and for each example we compute the mean log-probability over reasoning paths of the same type. Empirically, we find that **valid reasoning paths lead to substantially higher ground-truth answer probabilities than invalid ones**, even though most of the invalid reasoning paths in our analysis also end with the correct answer. In other words, reasoning paths that contain hallucinations or reward hacking are naturally assigned lower reward under group-relative reward normalization and are suppressed during training.

Specifically, the analysis is as follows. For each example and each reasoning type that is present on that example, we compute the improvement

$$\Delta \log p = \log p(y^* | \text{reasoning type}) - \mathbb{E}[\log p(y^* | \text{question-only})], \quad (23)$$

where $\mathbb{E}[\log p(y^* | \text{question-only})]$ is the mean log-probability of the ground-truth answer y^* over the eight question-only reasoning paths. Table 2 summarizes the distribution of $\Delta \log p$ across exam-

Reasoning type	n	mean	std	min	25%	50%	75%	max
Valid & Independent	53	4.97	4.25	-7.18	3.63	5.45	7.41	12.38
Valid & Dependent	31	4.11	4.36	-6.23	0.54	5.06	7.29	11.03
Invalid & Independent	62	2.09	3.93	-7.71	-0.72	1.20	5.16	10.80
Invalid & Dependent	59	3.27	4.15	-6.38	-0.04	3.97	5.95	13.65

Table 2: Per-example improvement $\Delta \log p$ for answer-conditioned reasoning types. Here n is the number of examples that contain at least one reasoning path of the given type; statistics are computed over the per-example means. Valid and independent reasoning achieves the largest mean and median gains, with a notably high 25th percentile (3.63), indicating strong and stable improvements over the question-only baseline and other reasoning types at the distribution level.

Comparison	n	t	p
Valid & Independent vs. Valid & Dependent	28	-1.54	0.134
Valid & Independent vs. Invalid & Independent	34	2.93	0.006
Valid & Independent vs. Invalid & Dependent	32	2.15	0.039

Table 3: Paired t -tests on per-example $\Delta \log p$, comparing valid, independent reasoning with other answer-conditioned types. Each test is computed on examples where both types are present.

ples for each reasoning type. All four types yield positive mean improvements over the question-only baseline, but valid and independent reasoning achieves the largest mean and median gains, with a particularly high 25th percentile (3.63), indicating both strong and stable improvements.

To compare reasoning types on the same examples, we run paired t -tests between valid and independent reasoning and each of the other three types, restricted to examples where both types are present (Table 3). Valid and independent reasoning significantly outperforms both invalid types in terms of $\Delta \log p$ ($p < 0.05$), while the difference between valid and independent reasoning and valid but dependent reasoning is not statistically significant.

To conduct a fine-grained examination of whether answer-conditioned reasoning may induce shortcut learning or post-hoc rationalization, we conduct a rubric-based evaluation of reasoning quality. Using the same 85 challenging questions, we collect 8 question-only reasoning traces and 8 answer-conditioned reasoning traces for each question, resulting in 680 traces for each setting. We then use GPT-5.1 as a strict rubric-based evaluator to score each reasoning trace along four dimensions on a 1–5 scale: *omissions* (missing essential steps or facts; lower is better), *hallucinations* (incorrect or fabricated claims; lower is better), *posthoc* (reverse-engineered reasoning with unsupported leaps; lower is better), and *overall_quality* (overall reasoning quality and reliability; higher is better). The full evaluation prompt is below.

Prompt Used for Rubric-Based Evaluation

System:

You are a strict rubric-based evaluator.

User:

You are evaluating the quality of a reasoning excerpt for a multiple-choice question.

Use the question and the ground-truth answer to judge factual accuracy.

Scoring (1–5):

- omissions: missing essential steps or key facts needed to justify the answer (1=none, 5=severe).

- hallucinations: incorrect or fabricated facts or claims (1=none, 5=severe).

- posthoc: reasoning appears reverse-engineered from the answer, with unsupported leaps (1=none, 5=severe).

- overall_quality: overall reasoning quality and reliability (1=very poor, 5=excellent).

Return ONLY a valid JSON object with keys: omissions, hallucinations, posthoc, overall_quality, brief_justification.

The brief_justification must be 1–2 sentences.

Question: [QUESTION]

Ground Truth Answer: [GROUND_TRUTH]

Reasoning Excerpt: [REASONING_EXCERPT]

Table 4 reports the overall comparison between question-only and answer-conditioned reasoning, together with a breakdown of answer-conditioned

traces by the same four reasoning types defined above. Overall, answer-conditioned reasoning improves all four dimensions relative to question-only reasoning, reducing omissions (2.300 \rightarrow 2.128), hallucinations (3.253 \rightarrow 2.791), and post-hoc behavior (3.784 \rightarrow 3.166), while substantially improving overall quality (1.971 \rightarrow 2.751). Among all subgroups, *valid and independent* reasoning achieves the strongest results on every metric, with the lowest omissions (1.495), lowest hallucinations (1.890), lowest post-hoc score (2.014), and highest overall quality (3.867). By contrast, the weakest subgroup is *invalid and dependent*, which has the highest post-hoc score (4.229). This is precisely the behavior that our leakage-aware reward design is intended to suppress.

Interestingly, we also observe that post-hoc scores remain relatively high in the question-only setting. To better understand this phenomenon, we inspect individual completions rather than relying only on aggregate averages. We find that even without explicit access to the reference answer, the model can still exhibit post-hoc style behavior, such as unsupported factual add-ons, conclusion-first reasoning followed by retrofitted justification, or unstable switching between multiple incompatible lines of thought. At the completion level, these cases are usually not caused by answer leakage, but instead reflect *unstable reasoning* near the model’s capability boundary.

To make this pattern more precise, we further analyze the distribution of post-hoc scores. For question-only reasoning, posthoc = 4 accounts for 76.62% of traces and posthoc = 5 for 3.97%, yielding a combined posthoc ≥ 4 rate of 80.59%. For answer-conditioned reasoning, posthoc = 4 accounts for 29.56% and posthoc = 5 for 15.88%, yielding a combined posthoc ≥ 4 rate of 45.44%. This suggests that question-only traces are dominated by a broad accumulation of medium-high instability (many score-4 cases), whereas answer-conditioned traces exhibit a lower overall high-score rate but a sharper cluster of severe post-hoc failures (more score-5 cases), where local samples show stronger answer-alignment override behavior.

Within answer-conditioned reasoning, subgroup analysis further clarifies where severe post-hoc behavior is concentrated. As shown in Table 5, the strongest concentration occurs in the *invalid and dependent* subgroup, where 83.58% of traces have posthoc ≥ 4 , including 44.28% with the maximum score of 5. By contrast, the *valid and independent*

Reasoning type	n	omissions ↓	hallucinations ↓	posthoc ↓	overall_quality ↑
Question-only (overall)	680	2.300	3.253	3.784	1.971
Answer-conditioned (overall)	680	2.128	2.791	3.166	2.751
Valid & Independent	218	1.495	1.890	2.014	3.867
Valid & Dependent	55	1.873	2.182	3.000	3.145
Invalid & Independent	206	2.510	3.272	3.393	2.218
Invalid & Dependent	201	2.493	3.443	4.229	1.980

Table 4: Rubric-based evaluation of reasoning quality. Lower is better for omissions, hallucinations, and posthoc, while higher is better for overall_quality. Answer-conditioned reasoning improves all dimensions over question-only reasoning, and valid and independent reasoning achieves the strongest performance across all four metrics.

is_valid_reasoning	is_independent	n	mean posthoc	posthoc = 4	posthoc = 5	posthoc ≥ 4
false	false	201	4.23	39.30%	44.28%	83.58%
false	true	206	3.39	47.09%	6.80%	53.88%
true	false	55	3.00	34.55%	3.64%	38.18%
true	true	218	2.01	2.75%	1.38%	4.13%

Table 5: Distribution of post-hoc scores for answer-conditioned reasoning traces across subgroups. Severe post-hoc behavior is concentrated in the invalid and dependent subgroup, while valid and independent reasoning remains consistently low on both posthoc = 4 and posthoc = 5.

subgroup remains very low on both posthoc = 4 and posthoc = 5, with only 4.13% of traces reaching posthoc ≥ 4. Therefore, although both settings contain post-hoc style behavior, the forms are different: question-only reasoning is primarily characterized by widespread instability, while answer-conditioned reasoning concentrates severe post-hoc failures in low-quality subsets.

We believe these findings are closely related to a key distinction between our method and prior answer-conditioned approaches. Previous methods typically use the answer as an exposed conditioning signal and directly optimize a mapping of the form *question + answer* → *reasoning steps*. Such a formulation is naturally susceptible to post-hoc rationalization, since the model can generate superficially plausible justifications after already anchoring on the answer. In contrast, our method uses a *why-style* prompting strategy that explicitly asks the model to reason about *why* the provided answer is correct, and to do so in a first-person think-aloud style as if solving the problem from scratch. We conjecture that this design makes answer-conditioned reasoning less prone to superficial answer imitation and more likely to elicit a genuinely explanatory reasoning process.

Overall, results suggest that shortcut learning is not the dominant effect in our setting. Although answer-conditioned reasoning can still produce severe post-hoc failures in some low-quality subsets, our reward design systematically favors logically

sound and answer-independent trajectories. The strongest gains in answer probability, as well as the best rubric-based reasoning quality, consistently come from valid and independent reasoning rather than from dependent or logically flawed traces.

A.5 ELBO Derivation with the Amortized Answer-conditioned Posterior

$\pi_\theta(z | x)$ is the prior over reasoning paths given x , and $\pi_\theta(y^* | x, z)$ is the likelihood of the reference answer y^* under path z . The original optimization objective is maximizing the utility under the reasoning distribution as

$$\log \mathcal{J}(\theta) = \log \mathbb{E}_{z \sim \pi_\theta(z|x)} [\pi_\theta(y^* | x, z)]. \quad (24)$$

Now we introduce an amortized posterior $\pi_\theta(z | x, y^*)$ to aid learning of $\pi_\theta(z | x)$. Starting from equation 24, multiply and divide inside the expectation by $\pi_\theta(z | x, y^*)$, and change the sampling distribution, we have:

$$\begin{aligned} \log \mathcal{J}(\theta) &= \log \mathbb{E}_{z \sim \pi_\theta(z|x)} \left[\pi_\theta(y^* | x, z) \cdot \frac{\pi_\theta(z | x, y^*)}{\pi_\theta(z | x, y^*)} \right] \\ &= \log \mathbb{E}_{z \sim \pi_\theta(z|x,y^*)} \left[\frac{\pi_\theta(y^* | x, z) \pi_\theta(z | x)}{\pi_\theta(z | x, y^*)} \right]. \end{aligned} \quad (25)$$

Applying Jensen’s inequality to the concavity of natural logarithm function $\log(\cdot)$,

$$\begin{aligned} \log \mathcal{J}(\theta) &\geq \mathbb{E}_{z \sim \pi_\theta(z|x, y^*)} \left[\log \pi_\theta(y^* | x, z) \right. \\ &\quad \left. + \log \pi_\theta(z | x) - \log \pi_\theta(z | x, y^*) \right] \end{aligned} \quad (27)$$

$$\begin{aligned} &= \mathbb{E}_{z \sim \pi_\theta(z|x, y^*)} \left[\log \pi_\theta(y^* | x, z) \right] \\ &\quad + \underbrace{\mathbb{E}_{z \sim \pi_\theta(z|x, y^*)} \left[\frac{\log \pi_\theta(z | x)}{\log \pi_\theta(z | x, y^*)} \right]}_{\text{forms a -KL term}} \end{aligned} \quad (28)$$

$$\begin{aligned} &= \mathbb{E}_{z \sim \pi_\theta(z|x, y^*)} \left[\log \pi_\theta(y^* | x, z) \right] \\ &\quad - \mathbb{D}_{\text{KL}}[\pi_\theta(z | x, y^*) \| \pi_\theta(z | x)]. \end{aligned} \quad (29)$$

Equation 29 is the Evidence Lower Bound (ELBO) on $\log \mathcal{J}(\theta)$ with the amortized answer-conditioned posterior $\pi_\theta(z|x, y^*)$. The KL estimator introduced by Schulman (2020) is as follows:

$$\begin{aligned} &\mathbb{D}_{\text{KL}}[\pi_\theta(z|x, y^*) \| \pi_\theta(z|x)] \\ &\approx \frac{\pi_\theta(z_t | x, z_{<t})}{\pi_\theta(o_{i,t} | x, y^*, \tilde{z}, z_{<t})} \\ &\quad - \log \frac{\pi_\theta(z_t | x, z_{<t})}{\pi_\theta(z_t | x, y^*, \tilde{z}, z_{<t})} - 1. \end{aligned} \quad (30)$$

A.6 Implementation Details

Each experiment is trained on 32 NVIDIA H100 GPUS. For **RAVR**, GRPO and DAPO, the learning rate for the policy model is 1e-6. In each rollout step, we sample eight responses per prompt for a batch of 32 prompts using a temperature of 1, and subsequently perform 2 policy updates on the collected responses with a batch size of 128. We adopt the clip-higher strategy and set the clip threshold as 0.8 and 1.27. For VeriFree and RLPR, we use the default setting in their official code. During evaluation, we follow the official recommended setting of Qwen3-1.7B, i.e., set Temperature=0.6, TopP=0.95 and TopK=20. In the evaluation, we use GPT-4.1-mini for correctness judgement. To reduce the evaluation variance, we report the final Avg4 for general reasoning multiple-choice tasks

and Avg4 for math reasoning tasks. The max generation length for training and evaluation is 8192. For DeepMath-103K, we randomly selected 5,000 samples for the training set. For MMLU-pro, we used the subset of the original data containing 1,000 randomly sampled samples (Yu et al., 2025b).

A.7 Additional Experiments for Generalization Analysis

To complement the generalization analysis, we include two additional experiments. First, we evaluate our method on a larger base model, Qwen3-8B, to further examine scalability. Second, we conduct an out-of-distribution (OOD) evaluation on ThinkBench (Huang et al., 2025) to assess performance under more challenging distribution shifts.

To further examine scalability, we additionally run experiments on Qwen3-8B. We set the maximum generation length to 2048 tokens, since stronger models often reach correct solutions with fewer reasoning tokens on problems of comparable difficulty, while this setting also keeps training time and compute cost manageable. We train for 300 RL steps and observe that the evaluation performance largely plateaus by this point. All other settings follow those used in our 1.7B/4B experiments.

Table 6 reports the results. RAVER consistently outperforms GRPO across all evaluated benchmarks. Notably, the improvement on Qwen3-8B is even larger than what we observe on smaller models, suggesting that **RAVR** may become more effective as the base model becomes stronger.

Another potential concern is whether the proposed method remains effective on more challenging or genuinely out-of-distribution reasoning tasks. To directly address this question, we evaluate **RAVR** on ThinkBench (Huang et al., 2025), a recent benchmark designed for dynamic OOD evaluation of LLM reasoning. ThinkBench constructs OOD variants using two families of data generation mechanisms: scenario-level semi-fact data generation and attack-level semi-fact data generation.

We use AIME-500, the largest subset in ThinkBench. It is extracted from AIME prob-

Table 6: Results on Qwen3-8B.

Training Set	Method	GPQA-Diamond (Avg@4)	AIME 24 (Avg@8)	AIME 25 (Avg@8)	AMC 23 (Avg@8)	Minerva (Avg@8)	Average
DeepMath	+ GRPO	26.27	12.89	7.42	44.38	55.88	29.37
DeepMath	+ RAVER (Ours)	33.93	19.58	16.25	62.50	59.93	38.44

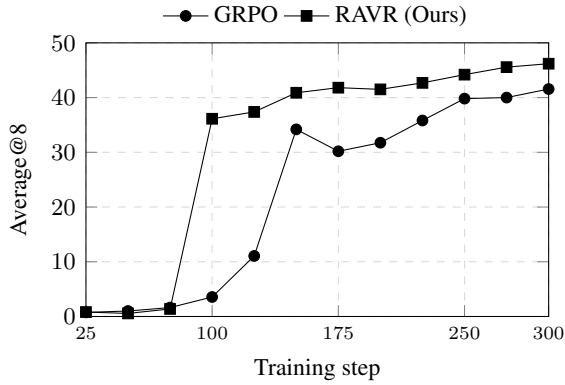


Figure 7: OOD evaluation on ThinkBench AIME-500 over training steps. RAVR converges faster, achieves a higher final performance, and exhibits a more stable training trajectory than GRPO.

lems from 1983 to 2023 and contains 500 original questions together with 2,000 OOD samples, i.e., four OOD variants per question. According to the original paper, nearly all models suffer noticeable degradation under OOD evaluation. For example, Qwen2.5-72B-IT drops from 43.2 under IID evaluation to 34.1 under OOD evaluation. The performance drop is typically more severe for weaker reasoning models. This makes ThinkBench a suitable testbed for evaluating whether our method improves robustness under distribution shift.

We split the 500 original questions into 450 training questions and 50 test questions, resulting in 1,800 training samples (450×4) and 200 test

samples (50×4). Because the ThinkBench paper mainly evaluates relatively large models, and models below 10B parameters generally achieve less than 10 accuracy on this benchmark, we conduct this experiment using Qwen3-8B. We set the maximum generation length to 2048 and train for 300 RL steps, evaluating on the test set every 25 steps. The metric is weighted average@8.

Figure 7 presents the results. RAVR outperforms GRPO throughout most of training and achieves a higher final score. More specifically, RAVR exhibits three advantages. First, it converges substantially faster: while GRPO only approaches its best performance around 250–300 steps, RAVR reaches a comparable level by around 150 steps. Second, it achieves a higher peak score, reaching 46.19 compared with 41.56 for GRPO. Third, it shows a more stable training trajectory: GRPO fluctuates noticeably between 150 and 225 steps, whereas RAVR improves much more steadily. These results suggest that the benefits of RAVR extend beyond in-distribution evaluation, and also transfer to challenging OOD reasoning settings.

A.8 More Reasoning Behavior Results

The reasoning behavior statistics on other datasets are reported in Figure 8.

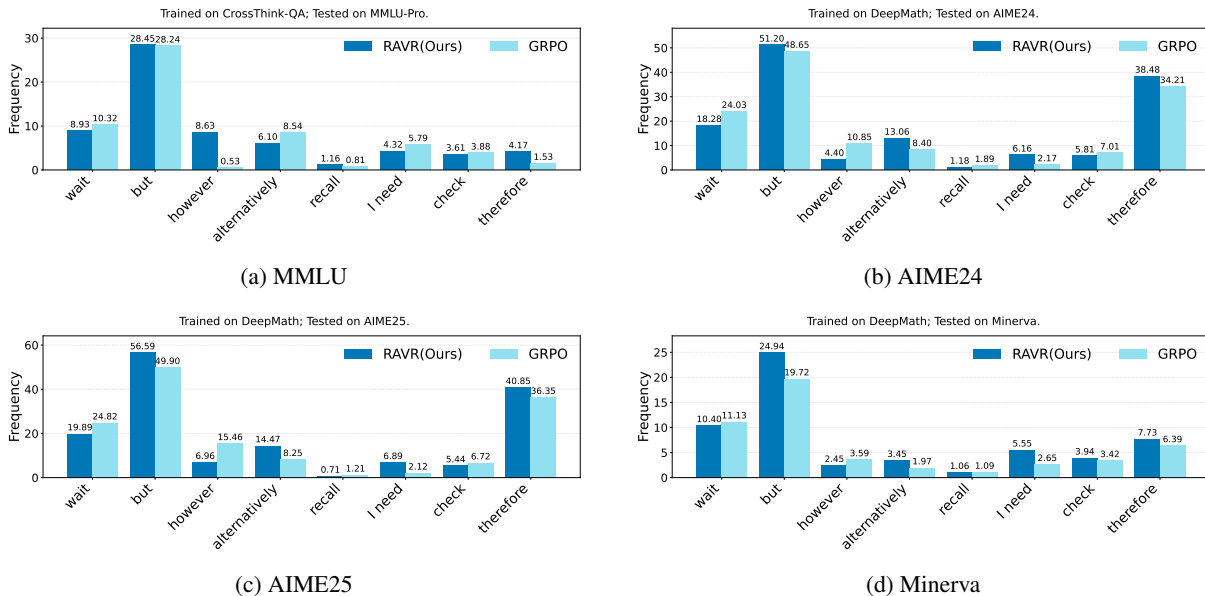


Figure 8: Reasoning behavior comparison across benchmarks. RAVR consistently shifts probability mass toward high-quality reasoning trajectories compared to GRPO, across both general (MMLU) and math-intensive benchmarks (AIME24, AIME25, Minerva).

A.9 Supplementary Learning Dynamics

As shown in Figure 9, we monitor whether the model’s responses contain explicit mentions of the phrase “reference answer” We find that, as training progresses, reasoning paths generated under reference-answer conditioning reveal the reference answer less frequently. Meanwhile, in the question-only setting, the model consistently maintains normal output behavior and never produces phrases such as “reference answer.”

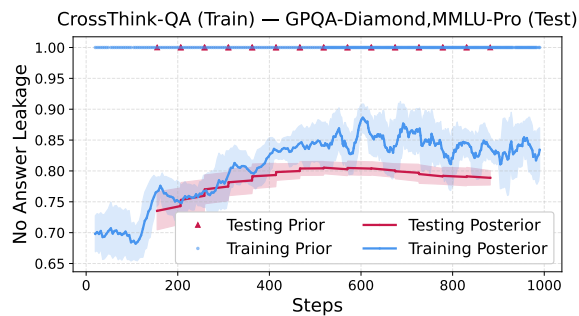


Figure 9: No Reference Answer Leakage