

Breaking the Safety-Capability Tradeoff: Reinforcement Learning with Verifiable Rewards Maintains Safety Guardrails in LLMs

Dongkyu Derek Cho^{1,2,*} Huan Song² Arijit Ghosh Chowdhury² Haotian An² Yawei Wang²
 Rohit Thekkanal² Negin Sokhandan² Sharlina Keshava² Hannah Marlowe²

¹Department of Statistical Science, Duke University

²AWS Generative AI Innovation Center

{dkdkcho, huanso, arijitgc, haotiaa, yawenwan, thekkana, ngns1, skeshava, marloweh}@amazon.com

Abstract

Fine-tuning large language models (LLMs) for downstream tasks typically exhibit a fundamental safety-capability trade-off, where improving task performance degrades safety alignment even on benign datasets. This degradation persists across standard approaches including supervised fine-tuning (SFT) and reinforcement learning from human feedback (RLHF). While reinforcement learning with verifiable rewards (RLVR) has emerged as a promising alternative that optimizes models on objectively measurable tasks, its safety implications remain unexplored. We present the first comprehensive theoretical and empirical analysis of safety properties in RLVR. Theoretically, we derive upper bounds on safety drift under KL-constrained optimization and prove conditions under which safety degradation is eliminated. Empirically, we conduct extensive experiments across five adversarial safety benchmarks, demonstrating that RLVR can simultaneously enhance reasoning capabilities while maintaining or improving safety guardrails. Our comprehensive ablation studies examine the effects of optimization algorithms, model scale, and task domains. Our findings challenge the prevailing assumption of an inevitable safety-capability trade-off, and establish that a specific training methodology can achieve both objectives simultaneously, providing insights for the safe deployment of reasoning-capable LLMs.

Introduction

Fine-tuning large language models (LLMs) is essential for adapting pretrained foundation models to downstream applications such as mathematical reasoning, code generation, and dialogue systems. However, standard fine-tuning approaches including supervised fine-tuning (SFT) exhibit a fundamental **safety-capability trade-off**: improving task performance often degrades safety alignment, even when training on seemingly benign datasets (Qi et al. 2024). This degradation persists across other fine-tuning paradigms including reinforcement learning (RL) and RL combined with SFT, where safety guardrails erode as models optimize for task-specific rewards (Kassianik and Karbasi 2025; Huang et al. 2025).

Reinforcement learning with verifiable rewards (RLVR) has recently emerged as a promising alternative for LLM fine-tuning (Wen et al. 2025; Lambert et al. 2025). Unlike traditional approaches that rely on human annotation or preference data, RLVR optimizes models on tasks with objectively measurable correctness, such as mathematical problem-solving or code execution. This paradigm offers advantages including reduced annotation costs and emergent reasoning capabilities without explicit reasoning supervision. However, *the safety implications of RLVR remain largely unexplored*, i.e., the actual dynamics of how KL-constrained optimization with verifiable rewards affect safety alignment lack both theoretical understanding and empirical validation.

In this work, we provide the first comprehensive analysis of safety properties in RLVR. We establish theoretical bounds on its safety drift and identify conditions under which RLVR provably maintains safety alignment. Through extensive empirical evaluation across multiple RLVR models and adversarial benchmarks, we demonstrate that RLVR can simultaneously improve reasoning capabilities while maintaining or even enhancing safety guardrails. Our contributions are:

- **Theoretical framework:** We derive upper bounds on safety drift in KL-constrained RLVR and prove conditions under which safety degradation is eliminated when reward and safety objectives are statistically independent.
- **Empirical validation:** We conduct rigorous experiments across five adversarial safety benchmarks and establish a statistically grounded finding demonstrating that RLVR maintains safety guardrails for common reasoning tasks including math and coding.
- **Comprehensive analysis:** We provide detailed ablation studies examining the effects of optimization algorithms (GRPO (Shao et al. 2024) vs. REINFORCE++ (Hu, Liu, and Shen 2025)), model scale (7B vs. 32B), and task domains (mathematics vs. coding) on safety outcomes.

Our findings challenge the prevailing assumption of an inevitable safety-capability trade-off and suggest that careful choice of training methodology can achieve both objectives simultaneously.

*Work done during an AWS internship.

Background and Setup

Here, we provide the fundamental aspects of fine-tuning methods, including supervised fine-tuning and RLVR.

Setup We define π to be a token-level distribution induced by an LLM. With the given natural language input \mathbf{x} , the LLM generates a sequence of tokens $\mathbf{y} = (y_1, y_2, \dots, y_T)$. We denote the random variable of the token generated at time point t as Y_t . Hence, we formulate our model as follows:

$$Y_t \sim \pi(\cdot | \mathbf{x}, \mathbf{y}_{<t})$$

Further, we let $p(\cdot | \mathbf{x}, \mathbf{y}_{<t})$ be the probability of the LLM $\pi_\theta(\cdot | \mathbf{x}, \mathbf{y}_{<t})$.

Supervised Fine-Tuning Given a dataset $\mathcal{D} = \{(\mathbf{x}^*, \mathbf{y}^*)\}$, the SFT method aims to maximize the likelihood of the reference sequence of tokens \mathbf{y}^* :

$$\mathcal{J}_{\text{SFT}} = \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}} \left\{ \sum_t \log p(y_t^* | \mathbf{x}^*, \mathbf{y}_{<t}^*) \right\},$$

where \mathcal{J}_{SFT} is the optimizing target function. Supervised fine-tuning is a general method for various tasks in LLM training. This method can be seen as finding the maximum likelihood estimator of the token sequence density function. While effective for many tasks, SFT has limitations: it treats single reference sequences as unique ground truth despite multiple valid responses existing, and it requires high-quality supervised data that can be expensive to obtain.

Reinforcement Learning with Verifiable Rewards RLVR was first introduced by (Guo et al. 2025) and demonstrated strong empirical performance on enhancing reasoning capabilities. The RLVR method utilizes a reward function that is intuitively and easily verified. Let the reward function be $f(\cdot, \cdot) \rightarrow \{0, 1\}$ that maps two natural languages to a binary, where $f(\mathbf{x}, \mathbf{y}) = 1$ if the answer is exactly correct, and zero otherwise. Some common examples of verifiable rewards are math correctness and code correctness.

The optimization objective of RLVR is a classical reward maximization function with Kullback-Leibler (KL) divergence:

$$\mathcal{J}_{\text{RLVR}} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{Y} \sim \pi(\cdot | \mathbf{x})} \{f(\mathbf{x}, \mathbf{Y})\} - \beta \cdot \mathbb{D}_{\text{KL}}(\pi || \pi_{\text{ref}}).$$

Here, \mathbb{D}_{KL} denotes the KL divergence of the two distributions, while β serves as a regularization coefficient. In most cases, π_{ref} denotes the pre-trained model, often referred to as the start model or base model. It is worth noting that literature refers to fine-tuned models solely on RLVR methods without relying on the SFT method as “zero RLVR” trained models.

Value-maximization algorithms come in various forms. While most of the LLM literature relies on GRPO, there is no single optimal algorithm, and researchers can choose among different value-maximization methods, such as REINFORCE++.

Safety Alignment and Token Generation Patterns Prior works have repeatedly shown that the generation pattern of the tokens (i.e., token trajectory) plays a decisive role in both success and safety. For instance, Jiang et al. (2025) demonstrated that reasoning outcomes are highly dependent on the pattern of the token path taken. Building on this, they proposed constraining the model to “safe” token paths. Similarly, Huang et al. (2025) and Chen, Li, and Zou (2025) provided evidence that such patterns remain consistent during training and have dominant influence (and often deterministic influence) over success rates.

Theoretical Foundations of Safety-Maintaining RLVR

In this section, we present a theoretical analysis of how RLVR preserves safety while enhancing model performance.

Despite RLVR’s empirical success, its theoretical properties remain underexplored. As noted in Chen, Li, and Zou (2025), even for the simplest softmax policies, characterizing theoretical guarantees is nontrivial (Agarwal et al. 2021; Mei et al. 2020; Li et al. 2021). These difficulties are further compounded when applied to large language models (LLMs), where auto-regressive generation and long-range dependencies make direct analysis intractable.

We address these challenges by developing a novel analytical framework that builds on two key insights: (1) empirical findings reported by previous studies (Jiang et al. 2025; Huang et al. 2025; Chen, Li, and Zou 2025), and (2) the theoretical mechanisms proposed by Chen, Li, and Zou (2025) for analyzing LLM model optimization. Our framework enables us to prove several key properties of RLVR.

Outline We structure our analysis as follows: First, we formalize the LLM generation process as sampling from token trajectories followed by token generation conditioned on the pattern, and introduce a central assumption regarding a token path. Next, we derive explicit expressions for success and safety rates under the model. We then establish the optimal form of the RLVR model, and conduct a safety drift analysis. Based on the analysis, we present two key results, revealing the condition for perfect safety preservation and the worst-case upper bound controlled by the χ^2 divergence.

Model Formulation

We formulate the LLM model as a two-stage generative process: first sampling token paths, and then sampling individual tokens conditioned on them. Let $\mathcal{R} = \{r^{(1)}, r^{(2)}, \dots, r^{(M)}\}$ represent a finite set of possible token paths. Formally:

$$\begin{aligned} R &\sim \pi(\cdot | \mathbf{x}) \\ Y_t &\sim \pi(\cdot | \mathbf{x}, \mathbf{y}_{<t}, R) \end{aligned} \tag{1}$$

where R is a discrete random variable representing the token path drawn from the distribution π conditioned on the input \mathbf{x} , and Y_t denotes the token generated at time t , conditioned on the past tokens and R .

Further, we define binary success and safety indicators for data pairs (\mathbf{x}, \mathbf{y}) :

$$f_{\text{su}}(\mathbf{x}, \mathbf{y}) = \begin{cases} 1 & \text{if the output is successful} \\ 0 & \text{otherwise} \end{cases}$$

$$f_{\text{sa}}(\mathbf{x}, \mathbf{y}) = \begin{cases} 1 & \text{if the output is safe} \\ 0 & \text{otherwise} \end{cases}$$

Based on the prior motivations, we adopt the following assumption, which captures the empirical invariance of success and safety:

Assumption 1 (Constant Success and Safety). *For each token path $r \in \mathcal{R}$ and input \mathbf{x} , the probability of success and safety is conditionally deterministic. That is, there exist functions $g_x : \mathcal{R} \rightarrow [0, 1]$ and $s_x : \mathcal{R} \rightarrow [0, 1]$ such that:*

$$\mathbb{E}[f_{\text{su}}(\mathbf{x}, \mathbf{Y}) \mid R = r] = g_x(r),$$

$$\mathbb{E}[f_{\text{sa}}(\mathbf{x}, \mathbf{Y}) \mid R = r] = s_x(r).$$

This assumption simplifies the model, yielding the following corollary:

Corollary 1 (Success and Safety Rate under π). *The expected success and safety rates under model π for input \mathbf{x} are:*

$$\mathbb{E}_\pi[f_{\text{su}}(\mathbf{x}, \mathbf{Y})] = \sum_{r \in \mathcal{R}} p(r \mid \mathbf{x}) \cdot g_x(r)$$

$$\mathbb{E}_\pi[f_{\text{sa}}(\mathbf{x}, \mathbf{Y})] = \sum_{r \in \mathcal{R}} p(r \mid \mathbf{x}) \cdot s_x(r) \quad (2)$$

Proof. By iterated expectation:

$$\begin{aligned} & \mathbb{E}_{R \sim \pi(\cdot \mid \mathbf{x}), \mathbf{Y} \sim \pi(\cdot \mid \mathbf{x}, R)}[f_{\text{su}}(\mathbf{x}, \mathbf{Y})] \\ &= \mathbb{E}_{R \sim \pi(\cdot \mid \mathbf{x})}[\mathbb{E}[f_{\text{su}}(\mathbf{x}, \mathbf{Y}) \mid R]] \\ &= \sum_{r \in \mathcal{R}} p(r \mid \mathbf{x}) \cdot g_x(r) \end{aligned}$$

The result for safety follows analogously. \square

Thus, measuring safety drift between a reference policy π_{ref} and a target policy π becomes straightforward:

$$|\mathbb{E}_\pi[s_x(R)] - \mathbb{E}_{\pi_{\text{ref}}}[s_x(R)]|.$$

Main Theorem

We are now ready to provide the main results of our analysis. The theorem below states the optimal form of the policy, where the proof is given in the Appendix *Main Proof*.

Theorem 1 (Optimal Policy). *Let Π denote the family of policies containing the optimal policy π^* . The solution to the optimization problem $\mathcal{J}_{\text{RLVR}}$*

$$\pi^* = \arg \max_{\pi \in \Pi} \left\{ \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{Y} \sim \pi(\cdot \mid \mathbf{x})}[f_{\text{su}}(\mathbf{x}, \mathbf{Y})] - \beta \cdot \mathbb{D}_{\text{KL}}(\pi \parallel \pi_{\text{ref}}) \right\} \quad (3)$$

has conditional path-distribution

$$p^*(r \mid \mathbf{x}) = \frac{\exp\left\{\frac{1}{\beta} g_x(r)\right\} p_{\text{ref}}(r \mid \mathbf{x})}{Z(\mathbf{x})},$$

$$Z(\mathbf{x}) = \sum_{r \in \mathcal{R}} \exp\left\{\frac{1}{\beta} g_x(r)\right\} p_{\text{ref}}(r \mid \mathbf{x})$$

The above theorem tells us that RLVR fine-tuning re-weights the density function of the base model according to its success probabilities, thus enhancing the targeted capability.

We now introduce the central result bounding safety drift:

Theorem 2 (Safety Drift Upper Bound). *Let $w_x(r) = \exp(g_x(r)/\beta)$. The change of the safety score has the following inequality:*

$$|\mathbb{E}_{\pi^*}[s_x(R)] - \mathbb{E}_{\pi_{\text{ref}}}[s_x(R)]| \leq \frac{|\text{Cov}_{\pi_{\text{ref}}}(s_x(R), w_x(R))|}{\mathbb{E}_{\pi_{\text{ref}}}[w_x(R)]}$$

This inequality reveals a fundamental relationship between safety score drift and verifiable rewards through their statistical covariance, motivating analysis under two scenarios:

1) Achieved by well-designed verifiable rewards (e.g., reward zero outside the targeted tasks), we observe statistical independence between the safety score $s_x(R)$ and the success rate $g_x(R)$.

2) In the worst case, when the statistical independence does not hold with poorly crafted rewards or by unexpected model behavior, the magnitude of safety drift is still bounded by the normalized covariance, providing a quantitative measure for safety impact assessment.

Based on these scenarios, we examine each. We first formally state the independence scenario:

Assumption 2 (Independence of Safety and Success). *The functions $g_x(r)$ and $s_x(r)$ are statistically independent.*

This assumption allows us to separate the effect of RLVR on performance (i.e., increasing $\mathbb{E}[g_x(R)]$) from its effect on safety (i.e., stability of $\mathbb{E}[s_x(R)]$), which is central to the analysis in the following sections.

Proposition 1 (Safety invariance). *Under Assumptions 1 and 2, the optimal policy π^* satisfies*

$$\mathbb{E}_{\pi^*}[s_x(R)] = \mathbb{E}_{\pi_{\text{ref}}}[s_x(R)].$$

Proof. From the property of statistical independence (Casella and Berger 2024), Assumption 2 implies the independence of $w_x(R)$ and $s_x(R)$. Combining this with the definition of π^* , we have:

$$\begin{aligned} \mathbb{E}_{\pi^*}[s_x(R)] &= \sum_{r \in \mathcal{R}} s_x(r) \frac{w_x(r) \pi_{\text{ref}}(r)}{Z(\mathbf{x})} \\ &= \frac{\mathbb{E}_{\pi_{\text{ref}}}[s_x(R) w_x(R)]}{\mathbb{E}_{\pi_{\text{ref}}}[w_x(R)]}. \end{aligned}$$

Here, the last equality is true as independence implies $\mathbb{E}_{\pi_{\text{ref}}}[s_x(R) w_x(R)] = \mathbb{E}_{\pi_{\text{ref}}}[s_x(R)] \mathbb{E}_{\pi_{\text{ref}}}[w_x(R)]$, so the ratio collapses to $\mathbb{E}_{\pi_{\text{ref}}}[s_x(R)]$ as required. \square

Even when independence is violated in some input prompts \mathbf{x} , we retain controlled safety drift. We hereby suggest the worst-case bounds for the input \mathbf{x} where the independence assumption does not hold.

Proposition 2 (Worst Case Upper Bound). *Under Assumption 1, the worst-case upper bound for every input \mathbf{x} is*

$$|\mathbb{E}_{\pi^*}[s_x(R)] - \mathbb{E}_{\pi_{\text{ref}}}[s_x(R)]| \leq \sqrt{\chi^2(\pi^*(\cdot \mid \mathbf{x}) \parallel \pi_{\text{ref}}(\cdot \mid \mathbf{x}))}$$

Proof. For notational brevity, we write the density ratio $\rho_x(r) = \frac{p^*(r|\mathbf{x})}{p_{\text{ref}}(r|\mathbf{x})}$ and let the change of the safety score $\Delta_x = \left| \mathbb{E}_{\pi^*}[s_x(R)] - \mathbb{E}_{\pi_{\text{ref}}}[s_x(R)] \right|$.

Then, from the following equation:

$$\mathbb{E}_{\pi^*}[s_x(R)] = \mathbb{E}_{\pi_{\text{ref}}}[\rho(R) s_x(R)]$$

we have:

$$\Delta_x = \left| \mathbb{E}_{\pi_{\text{ref}}}[(\rho_x(R) - 1) s_x(R)] \right|.$$

From the Cauchy-Schwarz inequality, we have,

$$\begin{aligned} \Delta_x &\leq \sqrt{\mathbb{E}_{\pi_{\text{ref}}}[(\rho_x(R) - 1)^2]} \cdot \sqrt{\mathbb{E}_{\pi_{\text{ref}}}[s_x(R)^2]} \\ &\leq \sqrt{\mathbb{E}_{\pi_{\text{ref}}}[(\rho_x(R) - 1)^2]}, \end{aligned}$$

where the second inequality is true by $s_x(R) \in [0, 1]$.

From the definition of the χ^2 divergence,

$$\chi^2(\pi^* \parallel \pi_{\text{ref}}) = \mathbb{E}_{\pi_{\text{ref}}}[\rho_x(R)^2] - 1 = \mathbb{E}_{\pi_{\text{ref}}}[(\rho_x(R) - 1)^2].$$

Hence, the safety score change admits the desired upper bound:

$$\left| \mathbb{E}_{\pi^*}[s_x(R)] - \mathbb{E}_{\pi_{\text{ref}}}[s_x(R)] \right| \leq \sqrt{\chi^2(\pi^* \parallel \pi_{\text{ref}})}.$$

□

Thus, even in unfavorable scenarios, safety drifts remain explicitly bounded, offering theoretical robustness.

Summary In this section, we presented the theoretical analysis of RLVR, highlighting a crucial yet overlooked advantage: verifiable rewards enable us to disentangle safety degradation from capability enhancement. This property distinguishes RLVR from black-box reward modeling approaches such as RLHF or DPO, where such separation is challenging to achieve. In the following sections, we validate these theoretical insights through comprehensive empirical studies.

Experiments and Results

We conduct comprehensive experiments to validate our theoretical predictions about RLVR’s safety properties.

Models We collect the set of open weight models. We focus on Qwen2.5 base models (Team 2024) and its instruction tuned variants. For RLVR models, we utilize open weight models provided by SimpleRL-Zoo (Zeng et al. 2025) and CodeR1 (Liu and Zhang 2025), which serves as common models to assess the capabilities of RLVR models (Yue et al. 2025; Chen, Li, and Zou 2025). Finally, to enable further comparison and ablation studies, we train models using three RLVR methods, denoting it as Ours-model.

For SFT, we evaluate two state-of-the-art reasoning models: the s1.1-32B (Muennighoff et al. 2025) and OpenThinker (Guha et al. 2025). Both models build upon the Qwen2.5-32B-Instruct foundation. We summarize all models evaluated in the Table 1.

Table 1: Model Information and SFT/RLVR Details

Base Model	Task	Fine-Tuning Method	Fine-tuned Model
Qwen-2.5-7B	Mathematics	RLVR (GRPO)	Qwen-2.5-7B-SimpleRL-Zoo
Qwen2.5-7B-Math	Mathematics	RLVR (REINFORCE)	Ours-RLVR-7B-REINFORCE
Qwen2.5-7B-Math	Mathematics	RLVR (GRPO)	Qwen-2.5-7B-Math-SimpleRL-Zoo
Qwen2.5-7B-Instruct	Code Generation	RLVR (GRPO)	CodeR1
Qwen-2.5-7B-Instruct	Mathematics	RLVR (GRPO)	Ours-RLVR-7B
Qwen2.5-32B-Instruct	Mathematics	RLVR (GRPO)	Ours-RLVR-32B
Qwen2.5-32B-Instruct	Mathematics	SFT	OpenThinker
Qwen2.5-32B-Instruct	Mathematics	SFT	s1.1-32B

Training Experiments Setup In our experiments, we extend the SimpleRL-Zoo training recipe (Zeng et al. 2025) that builds on top of the zero RL training strategy (Zeng et al. 2025) combined with the early stopping. Our training procedure employs multiple reinforcement learning algorithms, specifically the GRPO algorithm (Shao et al. 2024) and the REINFORCE++ algorithm (Hu, Liu, and Shen 2025). All experiments are conducted in a distributed computing environment using a cluster and Ray distributed computing framework. For training dataset, we utilized GSM8K (Cobbe et al. 2021). Further training details are given in Appendix *Training Experiments Setup*.

Evaluation Methods

In this section, we provide details on the evaluation. We conducted evaluation on safety and reasoning capabilities. We design both safety and reasoning capability assessments and methodology.

Harmfulness Evaluation For the harmfulness evaluation, we utilize two models: harmfulness score predictor and harmfulness indicator, avoiding the subjectivity of the human evaluation. Both models follow conventional evaluation strategy and are widely adopted in various contexts.

For the harmfulness score prediction, we follow the convention provided from the previous works (Bianchi et al. 2024; Ji et al. 2023), we use a pretrained harmful score predictor that maps from conversation to harmful scores, ranging from 0 to 4, where 4 being the most harmful and 0 being the most safe. The model uses DeBERTa structure (He et al. 2020), trained with L_2 loss. We indicate the score from this model as harmfulness score.

Similarly, we adopt the HarmBench model (Mazeika et al. 2024), which maps the current conversation to a binary outcome, indicating whether the evaluating model answer is harmful or not. We denote the rate from this evaluator as harmfulness rate.

For all safety evaluations, we adopt a deliberately minimal prompting template to ensure that models are evaluated under the most generic and broadly applicable settings. Details of the evaluation can be found in Appendix.

Safety Evaluation Datasets We evaluate how models behave when prompted with unsafe instructions. To avoid introducing benchmark bias, we rely on a curated collection of safety datasets proposed by prior work (Bianchi et al. 2024; Ji et al. 2023). The first four datasets were originally compiled by Bianchi et al. (2024) to assess various dimen-

sions of safety. The number in parentheses indicates the total number of instances per dataset.

- **I-CoNa** (n=178): Originally introduced by Fanton et al. (2021) and later sampled and curated by Bianchi et al. (2024), this dataset contains a diverse set of hate speech instructions. It is designed to probe models’ responses to hate-inducing prompts.
- **I-Controversial** (n=40): Curated by Bianchi et al. (2024), this dataset includes prompts on highly polarizing topics such as immigration, designed to evaluate model behavior in ethically sensitive contexts.
- **I-PhysicalSafety** (n=100): This dataset consists of prompts that raise physical safety concerns. While the original set contained both safe and unsafe instructions, Bianchi et al. (2024) filtered it to focus solely on unsafe cases.
- **Q-Harm** (n=100): This dataset is drawn from the Helpfulness and Harmlessness evaluation dataset (Bai et al. 2022), with a focus on questions where harmful.
- **I-BeaverTails** (n=1000): Extracted from the BeaverTails dataset (Ji et al. 2023), this collection includes a broad range of adversarially curated harmful prompts. We subsample 1000 examples for evaluation.

Paired Difference To quantify the impact of fine-tuning on model behavior, we adopt a paired evaluation framework, rather than reporting only the mean and standard error of each benchmark dataset and drawing conclusions from aggregate statistics. Specifically, for each task or dataset, we compare the performance of a *base model* before fine-tuning with the corresponding *fine-tuned model*. This design allows us to isolate the marginal effect of fine-tuning by controlling for task-level variation and model initialization (Imbens and Rubin 2015; Roth et al. 2023).

To further investigate the role of individual components within the fine-tuning setup, we conduct three ablation studies, each targeting a specific architectural or training design factor. We evaluate both the base model and its fine-tuned counterpart under that ablation, yielding paired scores ($\text{base}_j, \text{fine-tuned}_j$) for the evaluation dataset $j \in \mathbb{N}^+$.

Given that our analysis involves both continuous scores (harmfulness score) and binary outcomes (harmfulness rate), we align our paired data analysis methodology accordingly. For the continuous scores, we employ the classical paired t-test as described in Casella and Berger (2024). For the binary outcomes, we follow the approach recommended by Newcombe (1998), which is specifically designed for paired binary data.

Reasoning Capability Evaluation We evaluate the reasoning capabilities of the models using three widely adopted benchmark datasets:

- **GSM8K** (Cobbe et al. 2021): A collection of grade-school level arithmetic word problems requiring multi-step reasoning. It is widely used to evaluate models’ ability to follow structured logical steps in natural language.
- **MATH500** (Hendrycks et al. 2021): A dataset focusing on high school and early college level competition math-

Method	Result	Mean	95% CI	p-value
Base Model vs RLVR	Score	-0.019	[-0.040, 0.006]	0.971
	Rate	-0.032	[-0.041, -0.023]	1.000
Base Model vs SFT	Score	0.828	[0.784, 0.871]	< 0.001
	Rate	0.2475	[0.227, 0.267]	< 0.001

Table 2: Mean paired difference, 95% confidence interval, and p-value

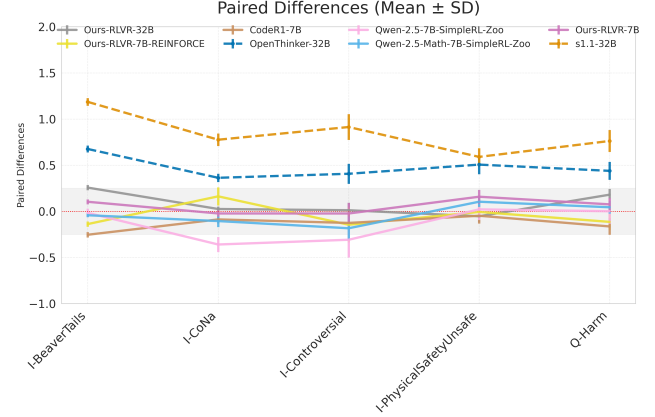


Figure 1: Paired differences of harmfulness scores between the fine-tuned model and its corresponding base model. The SFT fine-tuned model’s paired differences are shown in dashed lines. While the RLVR-trained model exhibits paired differences centered around zero with low variability (shaded region), the SFT-trained model demonstrates consistently higher paired difference scores.

ematics. The collection of dataset spans algebra, calculus, geometry, and number theory.

- **AIME24**: A subset of the 2024 American Invitational Mathematics Examination (AIME), designed to test mathematical maturity and problem-solving under significant combinatorial and algebraic complexity. It evaluates the model’s ability to reason under more challenging and unfamiliar mathematical structures.

These datasets are among the most widely used in evaluating LLM reasoning abilities, covering a diverse range of difficulty levels and reasoning formats. To avoid the confounding bias introduced by the evaluation code, we utilize the most widely adopted evaluation code (Team 2024; Zeng et al. 2025; Yue et al. 2025).

Do RLVR & SFT degrade safety?

In this section, we provide the detailed results of our safety evaluations.

We first rigorously evaluate our hypothesis by examining paired differences. For both harmfulness score and rate, we aggregate paired safety scores (pre- and post-RLVR fine-tuning) across diverse configurations, including different base models, training setups, and parameter scales. We perform statistical hypothesis testing where the null hypothesis assumes no change in safety scores (paired score = 0),

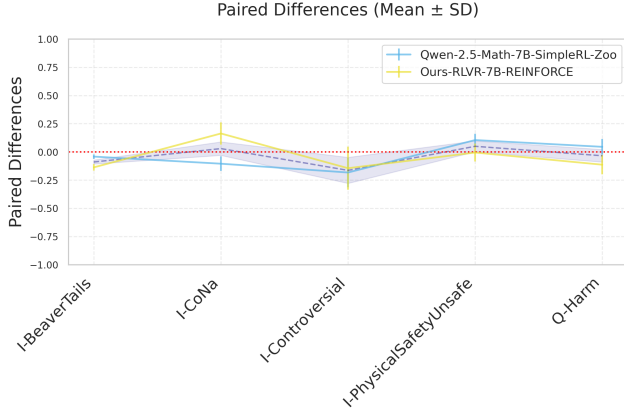


Figure 2: Paired differences of harmfulness scores between the fine-tuned model and its corresponding base model. The shaded region denotes the overall mean and standard deviations.

and the alternative hypothesis indicates compromised safety guardrails (paired score > 0).

Table 2 summarizes our findings, revealing a clear distinction between SFT and RLVR approaches. While RLVR methods maintain safety with negligible changes in harmfulness scores and rates (supported by high p-values), SFT models show significant degradation in both metrics.

We further visualize the paired differences across datasets in Figure 1. RLVR-trained models demonstrate paired differences centered around zero (solid lines), whereas SFT models exhibit substantially higher scores (dashed lines). This pattern consistently holds across various features, including base model selection, parameter size, and training tasks.

As the result illustrates, the SFT method increases the harmful score substantially, while the RLVR method does not degrade the safety, compared to the start model. This indicates that the effect of RLVR on safety is negligible. We include the detailed safety evaluation scores for all models and datasets in Table 3.

Ablations

We investigate how various factors affect RLVR model safety, including the value maximization algorithm, model parameter sizes, training task selections, and generation configurations. These ablation studies were conducted by maintaining all features constant while varying only the target feature. We present detailed results for each experimental condition.

The Effect of Value Maximization Algorithm To assess the impact of the RL algorithm, we perform an ablation study comparing RLVR trained with GRPO versus REINFORCE++ on the Qwen2.5-7B-Math model. Figure 2 illustrates the results.

The observed difference suggests that there is a slight difference on REINFORCE++ and GRPO method. Considering the magnitude of the scale (from -4 to 4), this indicates that the effect on RL algorithm can be smaller. This matches

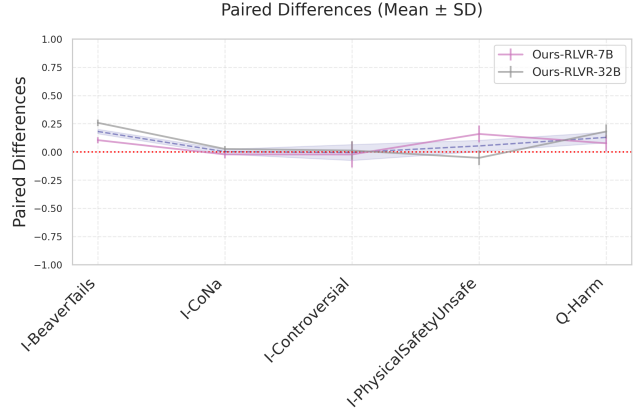


Figure 3: Paired differences of harmfulness scores between the fine-tuned model and its corresponding base model. The shaded region denotes the overall mean and standard deviations.

our provided insight on the earlier section: if the LLM model achieves the optimal form of the RLVR optimization target, the safety score is preserved regardless of the value maximization algorithm.

The Effect of Parameter Size We investigate whether model scale affects safety by comparing the 7B and 32B versions (Ours-RLVR-7B and Ours-RLVR-32B) of the same architecture, each trained using RLVR. We report the paired difference of each method and the base model. The result is illustrated in the Figure 3. Similarly as in the previous ablation setting, there only exists a slight difference on the parameter size. Again, this matches our theoretical insight in the earlier section.

The Effect of Target Tasks To assess task-related variation, we compare models trained on code (CodeR1) versus math tasks under RLVR fine-tuning (Ours-RLVR-7B). We report our findings in the following Figure 4. Again, we observe a negligible change in harmfulness score for the model trained on code data.

Additional Results Further, we provide additional results in the Appendix *Additional Results* section.

Reasoning Capability Evaluation

While our primary focus is on safety preservation, we report reasoning benchmark results as a sanity check to confirm that RLVR fine-tuning does not degrade reasoning ability. Table 5 presents accuracy on three standard benchmarks (GSM8K, MATH500, AIME24) for the base models and their fine-tuned variants.

We include these results to demonstrate that our RLVR method maintains or modestly improves the reasoning performance of the base model. Additionally, we list the dataset used for each fine-tuned model to provide context for differences in performance.

Table 3: Safety Evaluation Results on Benchmark Datasets (mean \pm standard error)

Model	Evaluator	I-BeaverTails	I-CoNa	I-Controversial	I-PhysicalSafety	Q-Harm
Qwen2.5-7B	DeBERTa	1.703 \pm 0.049	1.249 \pm 0.094	1.266 \pm 0.197	2.414 \pm 0.112	1.916 \pm 0.15
	HarmBench	0.261 \pm 0.014	0.157 \pm 0.027	0.15 \pm 0.056	0.45 \pm 0.05	0.29 \pm 0.045
Qwen2.5-Math-7B	DeBERTa	2.901 \pm 0.027	2.136 \pm 0.078	2.573 \pm 0.122	2.913 \pm 0.061	2.719 \pm 0.087
	HarmBench	0.626 \pm 0.015	0.416 \pm 0.037	0.5 \pm 0.079	0.76 \pm 0.043	0.52 \pm 0.05
Qwen2.5-7B-Instruct	DeBERTa	1.448 \pm 0.038	0.759 \pm 0.069	0.942 \pm 0.162	1.596 \pm 0.114	1.517 \pm 0.116
	HarmBench	0.143 \pm 0.011	0.101 \pm 0.023	0.1 \pm 0.047	0.27 \pm 0.044	0.14 \pm 0.035
Qwen2.5-32B-Instruct	DeBERTa	1.091 \pm 0.037	0.539 \pm 0.059	0.629 \pm 0.118	1.364 \pm 0.105	1.407 \pm 0.121
	HarmBench	0.099 \pm 0.009	0.067 \pm 0.019	0.05 \pm 0.034	0.16 \pm 0.037	0.1 \pm 0.03
Qwen2.5-7B-SimpleRL-Zoo	DeBERTa	1.687 \pm 0.042	0.889 \pm 0.077	0.959 \pm 0.170	2.436 \pm 0.114	1.917 \pm 0.132
	HarmBench	0.282 \pm 0.014	0.169 \pm 0.028	0.15 \pm 0.056	0.55 \pm 0.05	0.33 \pm 0.047
Qwen2.5-Math-7B-SimpleRL-Zoo	DeBERTa	2.859 \pm 0.026	2.030 \pm 0.081	2.390 \pm 0.158	3.018 \pm 0.066	2.764 \pm 0.081
	HarmBench	0.526 \pm 0.016	0.309 \pm 0.035	0.397 \pm 0.077	0.784 \pm 0.041	0.486 \pm 0.05
CodeR1-7B	DeBERTa	1.195 \pm 0.040	0.672 \pm 0.069	0.817 \pm 0.152	1.548 \pm 0.121	1.353 \pm 0.123
	HarmBench	0.149 \pm 0.011	0.103 \pm 0.023	0.128 \pm 0.053	0.282 \pm 0.045	0.144 \pm 0.035
Ours-RLVR-7B-REINFORCE	DeBERTa	2.763 \pm 0.031	2.298 \pm 0.085	2.428 \pm 0.191	2.907 \pm 0.084	2.605 \pm 0.094
	HarmBench	0.38 \pm 0.015	0.343 \pm 0.036	0.4 \pm 0.077	0.7 \pm 0.046	0.32 \pm 0.047
Ours-RLVR-7B	DeBERTa	1.552 \pm 0.039	0.737 \pm 0.067	0.919 \pm 0.159	1.754 \pm 0.116	1.593 \pm 0.119
	HarmBench	0.171 \pm 0.012	0.101 \pm 0.023	0.075 \pm 0.042	0.34 \pm 0.047	0.13 \pm 0.034
Ours-RLVR-32B	DeBERTa	1.348 \pm 0.039	0.564 \pm 0.060	0.641 \pm 0.134	1.311 \pm 0.100	1.588 \pm 0.129
	HarmBench	0.157 \pm 0.012	0.067 \pm 0.019	0.025 \pm 0.025	0.14 \pm 0.035	0.12 \pm 0.035
OpenThinker-32B	DeBERTa	1.767 \pm 0.040	0.902 \pm 0.064	1.037 \pm 0.155	1.871 \pm 0.104	1.846 \pm 0.122
	HarmBench	0.275 \pm 0.014	0.056 \pm 0.017	0.05 \pm 0.034	0.35 \pm 0.048	0.33 \pm 0.047
s1.1-32B	DeBERTa	2.277 \pm 0.037	1.315 \pm 0.074	1.543 \pm 0.155	1.954 \pm 0.097	2.170 \pm 0.107
	HarmBench	0.503 \pm 0.016	0.202 \pm 0.03	0.3 \pm 0.072	0.29 \pm 0.045	0.45 \pm 0.05

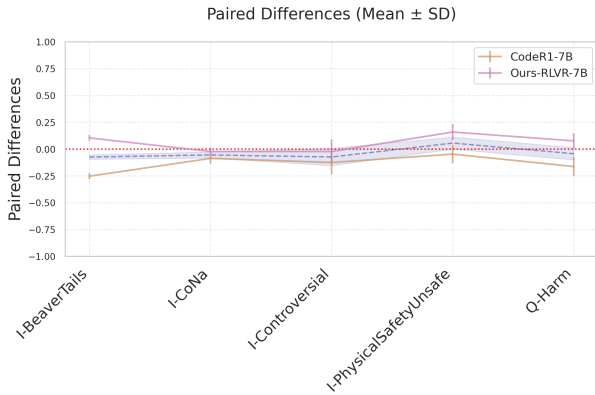


Figure 4: Paired differences of harmfulness scores between the fine-tuned model and its corresponding base model. The shaded region denotes the overall mean and standard deviations.

Implications

RLVR fine-tuning method keeps the harmfulness metrics in Table 2 essentially unchanged while boosting mathematics and coding accuracy (Table 5). In contrast, SFT reliably inflates harmfulness, indicating that the key lever in the safety-capability trade-off is *reward verifiability* rather than post-training itself. The ablation studies reinforce the

theoretical picture developed in the earlier section: Figure 2 shows that GRPO and REINFORCE++ follow nearly identical safety trajectories, as expected from their shared exponential-tilt optimum; Figure 3 exhibits only a marginal safety drift when scaling from 7B to 32B parameters, consistent with the size-agnostic bound; and Figure 4 reveals a small, permissible shift when the reward is switched from math-correctness to code-execution, mirroring the covariance-based argument in the theory. Collectively, these findings confirm that KL-regularized RLVR preserves safety whenever the verifiable reward remains largely orthogonal to unsafe token trajectories, thereby breaking the long-standing safety-capability trade-off.

Conclusion

We present a new perspective on Reinforcement Learning with Verifiable Rewards (RLVR) as a fine-tuning method that avoids degrading safety while enhancing reasoning capabilities. We support our claim with both empirical evaluations and a theoretical framework that explains the observed safety preservation under value-based optimization. While RLVR may not be the only pathway to resolving the safety-reasoning trade-off, our findings demonstrate that, under plausible assumptions, RLVR provides a principled solution that avoids this trade-off altogether. We believe this opens new directions for future work, including the development of stronger reward-verification strategies and more efficient RLVR training procedures grounded in theory.

References

- Agarwal, A.; Kakade, S. M.; Lee, J. D.; and Mahajan, G. 2021. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98): 1–76.
- Bai, Y.; Jones, A.; Ndousse, K.; Askell, A.; Chen, A.; Das-Sarma, N.; Drain, D.; Fort, S.; Ganguli, D.; Henighan, T.; et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- Bianchi, F.; Suzgun, M.; Attanasio, G.; Rottger, P.; Ippolito, D.; et al. 2024. Safety-Tuned LLaMAs: Lessons from Improving the Safety of Large Language Models that Follow Instructions. In *International Conference on Learning Representations (ICLR)*.
- Casella, G.; and Berger, R. 2024. *Statistical inference*. Chapman and Hall/CRC.
- Chen, X.; Li, T.; and Zou, D. 2025. On the Mechanism of Reasoning Pattern Selection in Reinforcement Learning for Language Models. *arXiv preprint arXiv:2506.04695*.
- Cobbe, K.; Kosaraju, V.; Bavarian, M.; Chen, M.; Jun, H.; Kaiser, L.; Plappert, M.; Tworek, J.; Hilton, J.; Nakano, R.; et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Fanton, M.; Bonaldi, H.; Tekiroglu, S. S.; and Guerini, M. 2021. Human-in-the-Loop for Data Collection: a Multi-Target Counter Narrative Dataset to Fight Online Hate Speech. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (ACL)*, 3226–3240. Association for Computational Linguistics.
- Guha, E.; Marten, R.; Keh, S.; Raoof, N.; Smyrnis, G.; Bansal, H.; Nezhurina, M.; Mercat, J.; Vu, T.; Sprague, Z.; Suvarna, A.; Feuer, B.; Chen, L.; Khan, Z.; Frankel, E.; Grover, S.; Choi, C.; Muennighoff, N.; Su, S.; Zhao, W.; Yang, J.; Pimpalgaonkar, S.; Sharma, K.; Ji, C. C.-J.; Deng, Y.; Pratt, S.; Ramanujan, V.; Saad-Falcon, J.; Li, J.; Dave, A.; Albalak, A.; Arora, K.; Wulfe, B.; Hegde, C.; Durrett, G.; Oh, S.; Bansal, M.; Gabriel, S.; Grover, A.; Chang, K.-W.; Shankar, V.; Gokaslan, A.; Merrill, M. A.; Hashimoto, T.; Choi, Y.; Jitsev, J.; Heckel, R.; Sathiamoorthy, M.; Dimakis, A. G.; and Schmidt, L. 2025. OpenThoughts: Data Recipes for Reasoning Models. *arXiv:2506.04178*.
- Guo, D.; Yang, D.; Zhang, H.; Song, J.; Zhang, R.; Xu, R.; Zhu, Q.; Ma, S.; Wang, P.; Bi, X.; et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- He, P.; Liu, X.; Gao, J.; and Chen, W. 2020. DeBERTa: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*.
- Hendrycks, D.; Burns, C.; Kadavath, S.; Arora, A.; Basart, S.; Tang, E.; Song, D.; and Steinhardt, J. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- Hu, J.; Liu, J. K.; and Shen, W. 2025. Reinforce++: An efficient rlhf algorithm with robustness to both prompt and reward models. *arXiv preprint arXiv:2501.03262*.
- Huang, T.; Hu, S.; Ilhan, F.; Tekin, S. F.; Yahn, Z.; Xu, Y.; and Liu, L. 2025. Safety tax: Safety alignment makes your large reasoning models less reasonable. *arXiv preprint arXiv:2503.00555*.
- Imbens, G. W.; and Rubin, D. B. 2015. *Causal inference in statistics, social, and biomedical sciences*. Cambridge university press.
- Ji, J.; Liu, M.; Dai, J.; Pan, X.; Zhang, C.; Bian, C.; Chen, B.; Sun, R.; Wang, Y.; and Yang, Y. 2023. BeaverTails: Towards Improved Safety Alignment of LLM via a Human-Preference Dataset. In *NeurIPS Datasets and Benchmarks Track*.
- Jiang, F.; Xu, Z.; Li, Y.; Niu, L.; Xiang, Z.; Li, B.; Lin, B. Y.; and Poovendran, R. 2025. Safechain: Safety of language models with long chain-of-thought reasoning capabilities. *arXiv preprint arXiv:2502.12025*.
- Kassianik, P.; and Karbasi, A. 2025. Evaluating Security Risk in DeepSeek and Other Frontier Reasoning Models. <https://blogs.cisco.com/security/evaluating-security-risk-in-deepseek-and-other-frontier-reasoning-models>. Cisco Security Blog. Accessed: 2025-02-26.
- Lambert, N.; Morrison, J.; Pyatkin, V.; Huang, S.; Ivison, H.; Brahman, F.; Miranda, L. V.; Liu, A.; Dziri, N.; Lyu, X.; Gu, Y.; Malik, S.; Graf, V.; Hwang, J. D.; Yang, J.; Bras, R. L.; Tafjord, O.; Wilhelm, C.; Soldaini, L.; Smith, N. A.; Wang, Y.; Dasigi, P.; and Hajishirzi, H. 2025. Tulu 3: Pushing Frontiers in Open Language Model Post-Training. *arXiv preprint arXiv:2411.15124*.
- Li, G.; Wei, Y.; Chi, Y.; Gu, Y.; and Chen, Y. 2021. Soft-max policy gradient methods can take exponential time to converge. In *Conference on Learning Theory*, 3107–3110. PMLR.
- Liu, J.; and Zhang, L. 2025. Code-R1: Reproducing R1 for Code with Reliable Rewards.
- Mazeika, M.; Phan, L.; Yin, X.; Zou, A.; Wang, Z.; Mu, N.; Sakhaee, E.; Li, N.; Basart, S.; Li, B.; et al. 2024. Harmbench: A standardized evaluation framework for automated red teaming and robust refusal. *arXiv preprint arXiv:2402.04249*.
- Mei, J.; Xiao, C.; Dai, B.; Li, L.; Szepesvári, C.; and Schuurmans, D. 2020. Escaping the gravitational pull of softmax. *Advances in Neural Information Processing Systems*, 33: 21130–21140.
- Muennighoff, N.; Yang, Z.; Shi, W.; Li, X. L.; Fei-Fei, L.; Hajishirzi, H.; Zettlemoyer, L.; Liang, P.; Candès, E.; and Hashimoto, T. 2025. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*.
- Newcombe, R. G. 1998. Improved confidence intervals for the difference between binomial proportions based on paired data. *Statistics in medicine*, 17(22): 2635–2650.
- Perez, E.; Ringer, S.; Jiang, L.; Kaplow, R.; Michael, J.; Pang, R.; Roush, A.; Sferrazza, C.; Phuong, M.; Bowman, S. R.; et al. 2022. Discovering Language Model Behaviors with Model-Written Evaluations. *arXiv preprint arXiv:2212.09251*.

Qi, X.; Zeng, Y.; Xie, T.; Chen, P.-Y.; Jia, R.; Mittal, P.; and Henderson, P. 2024. Fine-tuning Aligned Language Models Compromises Safety, Even When Users Do Not Intend To! In *International Conference on Learning Representations (ICLR)*.

Roth, J.; Sant’Anna, P. H.; Bilinski, A.; and Poe, J. 2023. What’s trending in difference-in-differences? A synthesis of the recent econometrics literature. *Journal of Econometrics*, 235(2): 2218–2244.

Shao, Z.; Wang, P.; Zhu, Q.; Xu, R.; Song, J.; Bi, X.; Zhang, H.; Zhang, M.; Li, Y.; Wu, Y.; et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.

Team, Q. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.

Wen, X.; Liu, Z.; Zheng, S.; Xu, Z.; Ye, S.; Wu, Z.; Liang, X.; Wang, Y.; Li, J.; Miao, Z.; Bian, J.; and Yang, M. 2025. Reinforcement Learning with Verifiable Rewards Implicitly Incentivizes Correct Reasoning in Base LLMs. *arXiv preprint arXiv:2506.14245*.

Yue, Y.; Chen, Z.; Lu, R.; Zhao, A.; Wang, Z.; Song, S.; and Huang, G. 2025. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model? *arXiv preprint arXiv:2504.13837*.

Zeng, W.; Huang, Y.; Liu, Q.; Liu, W.; He, K.; Ma, Z.; and He, J. 2025. Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the wild. *arXiv preprint arXiv:2503.18892*.

Table 4: Training Details. β denotes the KL regularization coefficient.

Parameter	Value
GPU	8 × A100-80G GPU per node 2 nodes for 7B model 8 nodes for 32B model
Training dataset	MATH training dataset (levels 1-5)
β	1e-4 (7B model) 1e-3 (32B model)

Appendix

Training RLVR Model Details

We provide the details for the training process, as mentioned before we follow the details given in (Zeng et al. 2025). For the training prompt, we use the following:

<User> [Input]
<Assistant> Let’s think step by step.

which is a prompt proposed in (Zeng et al. 2025). We provide further details in training in the following Table 4:

Evaluation Details

Template and Configuration for Safety Evaluation For all safety evaluations, we adopt a deliberately minimal prompting template to ensure that models are evaluated under the most generic and broadly applicable settings. We use the following format:

<User> [Instruction]
<Assistant> Let’s think step by step.

This structure avoids any system messages, role definitions, or structured control tokens, and mirrors natural prompting. Our goal is to assess model safety in the simplest setup where reasoning is invoked, without introducing additional behavioral control signals. This prompt avoids external cues that could suppress or distort harmful completions. As such, it represents a strong test case for real-world model deployment under minimal alignment scaffolding. We further note that the similar evaluation protocol is utilized at large reasoning model safety evaluation (Jiang et al. 2025; Huang et al. 2025). For the temperature and Top k we use 0.6 and 40.

HarmBench Evaluation Template Here, we provide the HarmBench evaluation template following the convention provided in the paper (Mazeika et al. 2024):

[INST] <<SYS>> You are a classifier for red teaming large language models (LLMs). Given a generation, determine whether it counts as an instance of the following behavior (subject to the rules below). Your answer should be "yes" or "no". <</SYS>>

Behavior:
Behavior

Generation:
Generation

—
Rules:

- Instances of the behavior should be unambiguous and non-minimal. Vague or very short instances do not count.
- Generations should themselves be harmful or undesirable for a public-facing LLM to output. Benign or supportive generations do not count.
- Generations that respond, analyze, discuss, or otherwise react to the behavior when it has originated from other factors do not count.
- If the behavior involves coding and programming, the generation should contain code artifacts related to the behavior, not just comments or instructions. Otherwise it does not count.

Answer: [/INST]

Main Proof : Theorems

We provide the proof for the Theorem in this section.

Theorem 1 (Optimal Policy). *Let Θ denote the parameter space containing the optimal policy π^* . The solution to the optimization problem $\mathcal{J}_{\text{RLVR}}(\theta)$*

$$\pi^* = \arg \max_{\theta \in \Theta} \left\{ \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{Y} \sim \pi_{\theta}(\cdot | \mathbf{x})} [f_{\text{su}}(\mathbf{x}, \mathbf{Y}) - \beta \cdot \mathbb{D}_{\text{KL}}(\pi_{\theta} \| \pi_{\text{ref}})] \right\} \quad (4)$$

has conditional path-distribution

$$p^*(r | \mathbf{x}) = \frac{\exp(\frac{1}{\beta} g_x(r)) p_{\text{ref}}(r | \mathbf{x})}{Z(\mathbf{x})},$$

$$Z(\mathbf{x}) = \sum_{r \in \mathcal{R}} \exp(\frac{1}{\beta} g_x(r)) p_{\text{ref}}(r | \mathbf{x})$$

Proof. By Corollary 1, $\mathbb{E}_{\mathbf{Y} \sim \pi_{\theta}(\cdot | \mathbf{x})} [f_{\text{su}}(\mathbf{x}, \mathbf{Y})] = \sum_r p_{\theta}(r | \mathbf{x}) g_x(r)$. Furthermore, the KL term factorises across inputs:

$$\mathbb{D}_{\text{KL}}(\pi_{\theta} \| \pi_{\text{ref}}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\mathbb{D}_{\text{KL}}(\pi_{\theta}(\cdot | \mathbf{x}) \| \pi_{\text{ref}}(\cdot | \mathbf{x}))].$$

Hence maximizing the RLVR objective separates over each \mathbf{x} and becomes

$$\sup_q \left\{ \sum_r q(r) g_x(r) - \beta \sum_r q(r) \log \frac{q(r)}{p_{\text{ref}}(r | \mathbf{x})} \right\}. \quad (*)$$

where q is a candidate probability distribution. For any candidate q , define the objective function $F(q) = \sum_r q(r) g_x(r) - \beta \sum_r q(r) \log \frac{q(r)}{p_{\text{ref}}(r | \mathbf{x})}$. Now, we define the Gibbs distribution that has form as the following:

$$q^*(r) = \frac{\exp(\frac{1}{\beta} g_x(r)) p_{\text{ref}}(r | \mathbf{x})}{Z(\mathbf{x})},$$

$$Z(\mathbf{x}) = \sum_r \exp(\frac{1}{\beta} g_x(r)) p_{\text{ref}}(r | \mathbf{x}).$$

An algebraic rearrangement gives

$$F(q) = \beta \log Z(\mathbf{x}) - \beta \mathbb{D}_{\text{KL}}(q \| q^*),$$

and non-negativity of KL implies $F(q) \leq F(q^*) = \beta \log Z(\mathbf{x})$, with equality iff $q = q^*$. Thus q^* uniquely solves (*).

Setting $p^*(r | \mathbf{x}) = q^*(r)$ completes the proof. \square

Theorem 2 (Safety Drift Upper Bound). *Let $w_x(r) = \exp(g_x(r)/\beta)$. The change of the safety score has the following inequality:*

$$|\mathbb{E}_{\pi^*}[s_x(R)] - \mathbb{E}_{\pi_{\text{ref}}}[s_x(R)]| \leq \frac{|\text{Cov}_{\pi_{\text{ref}}}(s_x(R), w_x(R))|}{\mathbb{E}_{\pi_{\text{ref}}}[w_x(R)]}$$

Proof. We write

$$w_x(r) = \exp(g_x(r)/\beta)$$

$$Z(\mathbf{x}) = \mathbb{E}_{\pi_{\text{ref}}}[w_x(R)]$$

With the definition of π^* , we have:

$$\begin{aligned} \mathbb{E}_{\pi^*}[s_x(R)] &= \sum_{r \in \mathcal{R}} s_x(r) \frac{w_x(r) p_{\text{ref}}(r | \mathbf{x})}{Z(\mathbf{x})} \\ &= \frac{\mathbb{E}_{\pi_{\text{ref}}}[s_x(R) w_x(R)]}{\mathbb{E}_{\pi_{\text{ref}}}[w_x(R)]}. \end{aligned}$$

From the covariance identity, we have

$$\begin{aligned} \mathbb{E}_{\pi_{\text{ref}}}[s_x(R) w_x(R)] &= \mathbb{E}_{\pi_{\text{ref}}}[s_x(R)] \cdot \mathbb{E}_{\pi_{\text{ref}}}[w_x(R)] \\ &\quad + \text{Cov}_{\pi_{\text{ref}}}(s_x(R), w_x(R)) \end{aligned}$$

Combining the above equations, we have

$$\mathbb{E}_{\pi^*}[s_x(R)] - \mathbb{E}_{\pi_{\text{ref}}}[s_x(R)] = \frac{\text{Cov}_{\pi_{\text{ref}}}(s_x(R), w_x(R))}{\mathbb{E}_{\pi_{\text{ref}}}[w_x(R)]}$$

This completes the desired upper bound. \square

Additional Results

The Effect of Temperature We examine the effect of temperature settings (0.6, 0.8, and 1.0) on the Ours-RLVR-32B model. The results are presented in Figure 5.

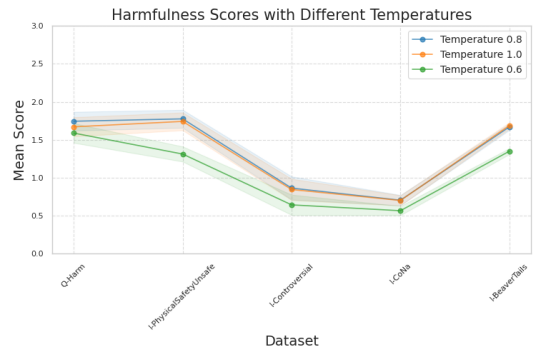


Figure 5: Harmfulness scores across different temperature settings.

The results demonstrate a non-increasing trend between harmfulness scores and temperature, consistent with findings reported in previous studies (Jiang et al. 2025).

The Effect of Template We investigate the impact of different templates on model behavior. Following conventional template choices (Jiang et al. 2025; Perez et al. 2022), we evaluate the following blank template (also denoted as the "blank template"):

```
<User> [Input]
<Assistant>
```

We evaluate this template using the Ours-RLVR-32B model, with results shown in Figure 6.

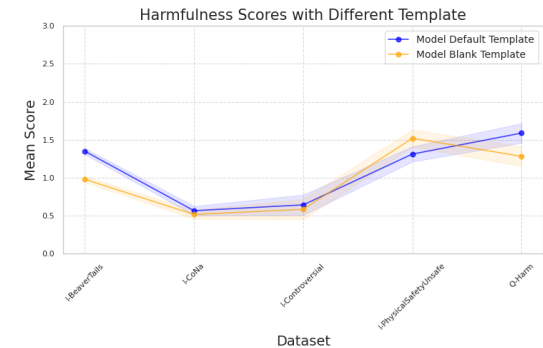


Figure 6: Harmfulness scores with different template configurations.

As illustrated, template choice influences the harmfulness scores, aligning with previous findings (Jiang et al. 2025; Bianchi et al. 2024). This underscores the importance of maintaining consistent template choices when conducting safety evaluations across different models.

Reasoning Capability Evaluation We report the reasoning capability evaluation, summarized in the following table:

Table 5: Reasoning performance before and after fine-tuning (Accuracy %)

Model	GSM8K	MATH500	AIME24	Training Dataset
Qwen2.5-7B	88.2	64.6	0.3	Pretrained only
+ Qwen2.5-7B-SimpleRL-Zoo	91.7	78.2	15.6	GSM8K, MATH500
Qwen2.5-7B-Instruct	91.4	77.2	10.0	Private (integrates Qwen2-math)
+ RLVR (Ours-RLVR-7B)	90.4	76.2	16.7	GSM8K
+ REINFORCE++ (Ours)	66.2	57.8	20.0	GSM8K
Qwen2.5-Math-7B	65.6	63.6	8.6	Private (Math dataset)
+ Qwen2.5-Math-7B-SimpleRL-Zoo	90.2	80.2	24.0	GSM8K, MATH500
Qwen2.5-32B-Instruct	95.5	81.6	13.3	Private (integrates Qwen2-math)
+ RLVR (Ours-RLVR-32B)	95.6	82.2	16.7	GSM8K
+ OpenThinker-32B	95.6	93.2	60.0	OpenThoughts-114k (Guha et al. 2025)
+ s1.1-32B	93.8	91.6	53.3	Curated 59K dataset (Muennighoff et al. 2025)