

JudgeRLVR: Judge First, Generate Second for Efficient Reasoning

Anonymous ACL submission

Abstract

Reinforcement Learning with Verifiable Rewards (RLVR) has become a standard paradigm for reasoning in Large Language Models. However, optimizing solely for final-answer correctness often drives models into aimless, verbose exploration, where they rely on exhaustive trial-and-error tactics rather than structured planning to reach solutions. While heuristic constraints like length penalties can reduce verbosity, they often truncate essential reasoning steps, creating a difficult trade-off between efficiency and verification. In this paper, we argue that discriminative capability is a prerequisite for efficient generation: by learning to distinguish valid solutions, a model can internalize a guidance signal that prunes the search space. We propose JudgeRLVR, a two-stage *judge-then-generate* paradigm. In the first stage, we train the model to judge solution responses with verifiable answers. In the second stage, we fine-tune the *same* model with vanilla generating RLVR initialized from the judge. Compared to Vanilla RLVR using the same math-domain training data, JudgeRLVR achieves a better quality–efficiency trade-off for Qwen3-30B-A3B: on **in-domain math**, it delivers about **+3.7 points** average accuracy gain with **-42%** average generation length; on **out-of-domain** benchmarks, it delivers about **+4.5 points** average accuracy improvement, demonstrating enhanced generalization.

1 Introduction

Reinforcement Learning with Verifiable Rewards (RLVR) has become a widely adopted method for improving large language models on complex reasoning tasks (Zhu et al., 2025; Shao et al., 2025). By providing sparse yet automatically checkable supervision signals (e.g., whether the final answer is correct), RLVR enables models to explore strategies beyond supervised demonstrations and has achieved strong results on reasoning benchmarks like math and code.

PROBLEM INPUT

Convert the point (0,3) in rectangular coordinates to polar coordinates. Enter your answer in the form (r, θ) , where $r > 0$ and $0 \leq \theta < 2\pi$.

Vanilla RLVR

Verbose

```
<think>
I need to .....

But here..... Let me
confirm..... I recall
that sometimes people
might think of  $\theta$  as
 $-\pi/2$  or  $3\pi/2$ , but
that would give
..... but here  $y$  is
positive.

Just to be
thorough..... But
the problem specifies
..... but here
.....

So, I'm confident.
</think>
```

⚠️ CoT tokens > 300

First, I need to.....The final answer is:

$(3, \pi/2)$

JudgeRLVR

Concise

```
<think>
We need to convert
(0,3) to polar
coordinates.  $r = \sqrt{x^2 + y^2} = 3$ .
Theta is angle from
positive x-axis to
(0,3). That's 90
degrees =  $\pi/2$  rad.
Since  $x=0, y>0$ , theta
=  $\pi/2$ . Final answer
 $\boxed{(3, \frac{\pi}{2})}$ .
</think>
```

⚡ CoT tokens < 100

First, I need to.....The final answer is:

$(3, \pi/2)$

Figure 1: Qualitative comparison of reasoning patterns between Vanilla RLVR and JudgeRLVR. While outcome-based RL guarantees final answer correctness, it often induces verbose and unstructured exploration (left). JudgeRLVR implicitly regularizes the reasoning process, leading to a direct and coherent solution path (right).

However, RLVR does not merely optimize correctness; it also implicitly shapes a model’s *thinking pattern*. When trained primarily with a final-answer reward, many LLMs drift toward a genera-

044
045
046
047

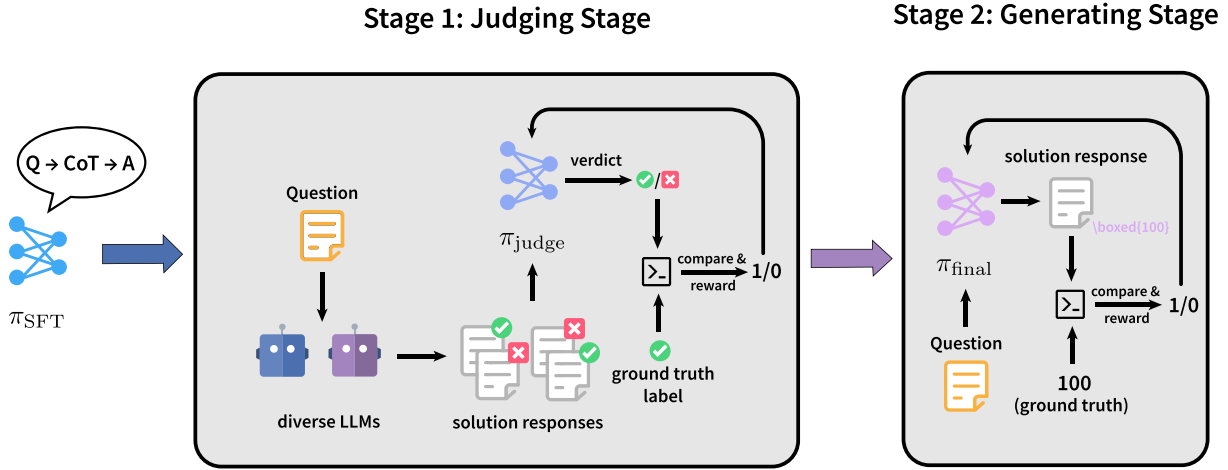


Figure 2: Pipeline for the two-stage training of JudgeRLVR

048 tive search style: enumerating tentative branches, 083
 049 revising intermediate steps, and performing explicit 084
 050 self-correction (Team et al., 2024). Figure 1 shows 085
 051 that the reasoning model’s output can be divided 086
 052 into 2 parts: chain-of-thought (CoT) trajectories 087
 053 and solution responses. CoT trajectories usually 088
 054 include all the details of the model’s lengthy thought 089
 055 process, whereas solution responses tend to present 090
 056 a clear logical process along with the final answer. 091
 057 The resulting CoT trajectories are often long and 092
 058 exhibit frequent backtracking (e.g., heavy use of 093
 059 “but”, “however”, “let’s try again”) (Cai and et al., 094
 060 2025). Long outputs are expensive, but a deeper 095
 061 issue is *quality*: such trajectories are low in infor- 096
 062 mation density. Existing methods such as Kimi 097
 063 K1.5 (Moonshot AI Team, 2025) and DAPO (Yu 098
 064 et al., 2025) introduce length penalties to reduce 099
 065 tokens and stabilize training, but often come with 100
 066 accuracy drops, making it difficult to achieve an 101
 067 ideal quality–efficiency trade-off. 102

068 This raises a more general question: can we 103
 069 explore a training paradigm that explicitly targets 104
 070 higher-quality thinking patterns, rather than hop- 105
 071 ing the model eventually “stumbles upon” them 106
 072 through repeated trial-and-error? Cognitive sci- 107
 073 ence suggests that expert reasoning is not defined 108
 074 by externalizing exhaustive search, but by *early 109*
 075 *discrimination and pruning*: unpromising paths 110
 076 are filtered before they are expanded (Chi et al., 111
 077 1981). This aligns with common intuition: experts 112
 078 quickly spot low-value reasoning paths and focus 113
 079 on high-yield hypotheses, while novices resort to 114
 080 trial-and-error and frequent backtracking. This in- 115
 081 tuition motivates the idea that, to achieve efficient 116
 082 reasoning, one should first develop the ability to 117

judge what constitutes direct and valid reasoning, and then learn how to reason for specific problems.

Based on the idea above, we propose **JudgeRLVR**, a two-stage *judge-then-generate* paradigm. As shown in Figure 2, in judging stage, we train the model as a judge: given candidate solution responses generated by diverse models, it learns to classify them as correct or incorrect. In generating stage, we initialize the model from the judging stage and apply *Vanilla RLVR* fine-tuning on the same model. Under this training paradigm, the model learns to implicitly prune low-quality branches prior to generating extensive textual search, thereby fostering more direct and reliable reasoning patterns. It is worth noting that, we do not impose any explicit length or quality penalties on the CoT trajectories; thus, any reduction in verbosity or refinement in reasoning style can be attributed to the discriminative priors transferred from the judging stage.

With our training paradigm, a model can achieve an optimal balance between performance and efficiency. Compared to Vanilla RLVR using the same training data, JudgeRLVR achieves a quality–efficiency trade-off on Qwen3-30B-A3B (Team, 2025): on **in-domain math**, it achieves about **+3.7 points** average accuracy improvement while reducing average generation length by about **42%**; on **out-of-domain** benchmarks, it delivers about **+4.5 points** average accuracy improvement. From extensive experiments, we also analyze the specific impact of this training paradigm on thinking patterns. We observe a shift in the linguistic style of the generated text, with notably fewer explicit backtracking cues such as “but” during reasoning. This

118 suggests that JudgeRLVR moves from an external-
119 ized trial-and-error process to an internal decision-
120 making process, thereby producing higher-quality
121 reasoning patterns.

122 In summary, our main contributions are: (i) we
123 propose **JudgeRLVR**, a new training paradigm for
124 inducing higher-quality thinking patterns via *judge-*
125 *then-generate*; (ii) JudgeRLVR demonstrates a bet-
126 ter quality–efficiency trade-off than Vanilla RLVR
127 on both in-domain and out-of-domain tasks; (iii)
128 we provide interpretable linguistic evidence which
129 is consistent with reduced explicit backtracking.

130 2 Related Work

131 2.1 RLVR and Verifiable-Reward Training for 132 Reasoning

133 Reinforcement Learning with Verifiable Rewards
134 (RLVR) is a leading approach for improving LLM
135 reasoning using automatic, task-specific rewards,
136 often outperforming supervised fine-tuning on
137 math and logic tasks (Lambert et al., 2024; Guo
138 et al., 2025a; Xu et al., 2025). However, RLVR
139 has notable limitations: it can yield gains even un-
140 der spurious rewards (Shao et al., 2025), and it
141 mainly reweights existing high-reward reasoning
142 paths rather than inducing new reasoning behav-
143 iors (Yue et al., 2025). This is consistent with
144 concerns that RLVR may implicitly encourage sub-
145 optimal thinking by optimizing final correctness in-
146 stead of process quality (DataLearner, 2025). Most
147 variants emphasize reward design or optimization,
148 leaving the quality–efficiency trade-off underex-
149 plored.

150 2.2 Mathematical Reasoning in Large 151 Language Models

152 Mathematical reasoning is a core benchmark for
153 LLMs, where recent progress relies heavily on
154 RLVR and targeted training (Team et al., 2024;
155 Shen et al., 2025). Data diversification via para-
156 phrasing can improve generalization (Shen et al.,
157 2025), while extended RLVR training improves
158 accuracy but tends to produce verbose reasoning
159 traces (Team et al., 2024). Synthetic data scaling
160 further boosts capability yet does not eliminate re-
161 dundancy (Zhou et al., 2025), and scaling studies
162 still report low information density in reasoning
163 outputs (Sun et al., 2024). Overall, improving both
164 correctness and efficiency remains challenging.

165 2.3 Learning to Judge: Evaluators, Verifiers, 166 and Discriminative Supervision

167 Training evaluators and using discriminative super-
168 vision can steer LLM outputs by separating good
169 from bad generations (Guo et al., 2025b; Tang et al.,
170 2025). Discriminative fine-tuning can suppress neg-
171 ative outputs without preference data (Guo et al.,
172 2025b), and stronger evaluators improve reasoning-
173 quality assessment via coherence-focused crite-
174 ria (Tang et al., 2025). Verifiers can detect rea-
175 soning errors with uncertainty signals (Han et al.,
176 2024), and systematic studies highlight the need
177 for task-specific validation metrics (Shimabucoro
178 et al., 2024). Yet existing work rarely adopts a two-
179 stage pipeline that first trains judgment and then
180 uses it to guide generation, which we address in
181 JudgeRLVR.

182 3 Method

183 **JudgeRLVR** divides the training of a single reason-
184 ing policy into two sequential stages with different
185 roles but intended capability transfer: a *judge* stage
186 that trains discriminative error awareness, followed
187 by a *generating* stage that optimizes solution gen-
188 eration. Our core hypothesis is that discriminative
189 competence is a prerequisite for efficient genera-
190 tion. Once the model internalizes what constitutes a
191 valid CoT trajectory, it can down-weight erroneous
192 branches early in generation without any explicit
193 length penalty.

194 Figure 2 illustrates the complete pipeline of our
195 two-stage training.

196 3.1 Notation and Task Definition

197 The dataset consists of problems in a specific do-
198 main: $x \in \mathcal{X}$ with gold final answers $a^*(x)$.

199 The **solution response** is denoted $z =$
200 (o_1, \dots, o_T) , a token sequence containing a clear
201 logical process and ending with a final answer
202 string. The final answer $a(z)$ is extracted from the
203 solution responses via a deterministic parser. We
204 do not train the model to judge directly on full out-
205 puts that include complete CoT trajectories because
206 such outputs are often extremely long and contain
207 a large amount of distracting information that can
208 interfere with the judgment. Instead, we expect the
209 model to judge directly based on the clean solution
210 response, so that it can identify errors with greater
211 precision.

212 We define a binary correctness label $\ell \in \{0, 1\}$:

$$213 \ell(x, z) = \mathbb{I}(a(z) = a^*(x)). \quad (1)$$

The policy (language model) is denoted π_θ .

3.2 Stage 1: Judging Stage

In judging stage, the model is trained as a judge. Given a problem x and a candidate solution response z , the policy produces a critique/commentary c which contains the CoT trajectory, and outputs a discrete verdict token $v \in \{0, 1\}$ (0 for incorrect and 1 for correct):

$$(c, v) \sim \pi_\theta(\cdot | x, z). \quad (2)$$

Data construction. We construct a discriminative dataset D_{judge} consisting of triplets (x, z, ℓ) :

- **Response generation (rollout).** For each problem x , we sample a set of candidate solution responses $\{z_1, \dots, z_n\}$ from several models, and obtain $\ell(x, z_i)$ by comparing $a(z_i)$ with $a^*(x)$.
- **Hard negative mining.** We prioritize *moderate-difficulty* problems whose empirical pass rates under rollouts is neither 0 nor 1, yielding more informative “nearly-correct” errors.
- **Class balancing.** For each x , we subsample to balance positives ($\ell = 1$) and negatives ($\ell = 0$), avoiding majority-class bias.

Reward. The reward is whether the final verdict v_i matches the true label $\ell(x, z_i)$:

$$r_i = \mathbb{I}(v_i = \ell(x, z_i)) \quad (3)$$

The training objective encourages v_i to match the true label $\ell(x, z_i)$. Meanwhile, the commentary c , as an explanatory prefix, is optimized by the same policy-gradient signal.

3.3 Stage 2: Generating Stage

In generating stage, the model is trained under Vanilla RLVR setting and outputs solutions of the problems. The policy is initialized from the judging stage weights.

Given a problem x , the policy generates a CoT trajectory and the solution response z (including the final parsed answer $a(z)$):

$$z \sim \pi_\theta(\cdot | x). \quad (4)$$

The reward remains the sparse binary correctness signal based only on the final answer:

$$r = \ell(x, z) = \mathbb{I}(a(z) = a^*(x)). \quad (5)$$

3.4 Hypothesized Mechanism and Verification Overview

We hypothesize that JudgeRLVR improves reasoning-pattern quality through the two stages in two corresponding ways: (i) **style transfer**: in judging stage, learning to judge induces a global shift in the model’s reasoning language style; (ii) **reduced backtracking**: in generating stage, training activates the internalized efficient pattern, reducing explicit self-correction and backtracking in text. We design targeted experiments to test both points; definitions and results are provided in Section 4.

4 Experimental Setup

4.1 Model Configuration and Training Data

We use Qwen3-30B-A3B (Team, 2025), a mixture-of-experts (MoE) architecture. To ensure a basic level of reasoning capability, we first conduct supervised fine-tuning (SFT) on curated open-source CoT datasets, obtaining Qwen3-30B-A3B-SFT with improved reasoning and instruction-following abilities.

For RLVR, we collect and filter open-source resources to build a dataset of 113k math problems with gold answers. To construct specialized training data for the judging stage, we adopt a strict sampling strategy:

- **Rollout generation.** For each problem, we use MiMo-7B RL (Xiaomi, 2025) and the target model Qwen3-30B-A3B-SFT to each generate 8 reasoning traces and answers, yielding 16 distinct reasoning paths.
- **Balanced sampling.** We filter outputs so that each problem retains an equal number of correct and incorrect solutions, avoiding bias from class imbalance.

Our training algorithm is DAPO (Yu et al., 2025), a GRPO-family policy gradient method. We implement dynamic sampling to filter out samples with pass rates 0 or 1 during rollout. The batch size is set to 256 and the rollout size is $n = 8$. We use AdamW optimizer (Loshchilov and Hutter, 2018) for training. The learning rate is set to 3×10^{-6} with a 10-step warmup. We use a training temperature of 1.0 and a testing temperature of 0.6, with top- p set to 1.0 and max tokens set to 65536 for both training and testing.

303	4.2 Evaluation Benchmarks	
304	We evaluate JudgeRLVR on in-domain math	351
305	reasoning: AIME24 (Zhang and Math-	352
306	AI, 2024), AIME25 (Zhang and Math-AI,	353
307	2025), MATH500 (Hendrycks et al., 2021),	354
308	HMMT_feb_2025 (Balunović et al., 2025),	355
309	BeyondAIME ([ByteDance-Seed], 2025) and	356
310	on diverse out-of-domain benchmarks: GPQA	357
311	Diamond (Science Reasoning) (Rein et al., 2023),	358
312	IFEval (Instruction Following) (Zhou et al., 2023),	359
313	LiveCodeBenchv6 (Coding) (Jain et al., 2024),	360
314	MMLU-Redux (General Knowledge) (Gema et al.,	361
315	2025), ZebraLogic (Logical Reasoning) (Lin et al.,	362
316	2025), covering a wide range of general capabili-	363
317	ties. The detailed setups of these benchmarks are	364
318	in appendix C.	365
319	4.3 Main Experiments	366
320	Our main experiments are designed to answer a	367
321	single overarching question: does the proposed	368
322	judge-then-generate paradigm improve the quality-	369
323	efficiency trade-off of RLVR beyond Vanilla	370
324	RLVR? To make this comparison meaningful, we	371
325	evaluate three training settings that isolate the ef-	372
326	fect of the staged paradigm from other factors.	373
327	Base SFT: the static baseline model Qwen3-	374
328	30B-A3B-SFT. Base SFT serves as a reference	375
329	point for the model’s reasoning ability before RL	376
330	and its default generation style, allowing us to mea-	377
331	sure how much of the downstream gains come	378
332	specifically from RLVR.	379
333	Vanilla RLVR: Vanilla generating RLVR opti-	380
334	mizing only for final-answer correctness. Vanilla	381
335	RLVR serves as a strong baseline to demonstrate	382
336	the advantages of our method. The total training is	383
337	250 steps.	384
338	JudgeRLVR: our judge-then-generate method.	385
339	The judging stage is 145 steps and the generating	386
340	stage is 105 steps.	387
341	Vanilla RLVR and JudgeRLVR use the same	388
342	training and evaluation hyperparameters and the	389
343	same total training steps. We report both accuracy	390
344	and generation length because our goal is not only	391
345	to increase correctness, but also to reduce unneces-	392
346	sary computation and improve information density	393
347	in the produced reasoning.	394
348	The generating and rollout prompt, and the judg-	395
349	ing prompt are in appendix A.	396
	4.4 Ablations	397
	We conduct targeted ablations to validate that	398
	JudgeRLVR’s gains come from the divided stages	399
	itself rather than from simply adding more training	
	signals.	
	Judge Only: only judging stage training. Judge	
	Only setting evaluates whether judge training alone	
	can already improve generation quality, or whether	
	generating stage is necessary to translate judging	
	competence into improved generation performance.	
	This ablation directly tests our core assumption	
	that judging is a prerequisite but not a substitute	
	for generating optimization.	
	Mixed Strategy: judging stage is replaced by	
	parallel training on mixed Judge/Generative data	
	(ratio 1:1), with each data type using its correspond-	
	ing reward; generating stage is then run as usual.	
	Mixed Strategy setting probes whether the order-	
	ing and separation of roles matters: by interleav-	
	ing judge-style and generation-style updates, we	
	test whether simultaneously optimizing two behav-	
	iors interferes with learning a clean internal deci-	
	sion policy. In other words, this ablation examines	
	whether the model benefits from first consolidating	
	an error-aware judge prior before being asked to	
	explore and exploit it during generation.	
	Together, these ablations aim to distinguish (i)	
	“benefit from judging” from (ii) “benefit from judg-	
	ing first, then generating,” clarifying which com-	
	ponent is responsible for improved efficiency and	
	generalization.	
	4.5 Mechanism Verification	
	To test style transfer and reduced backtracking, we	
	design:	
	Perplexity (PPL). PPL is the exponential of a	
	language model’s average negative log-likelihood	
	over a sequence; it measures how “surprised” or	
	“confused” the model is, and lower values indicate	
	the model is more familiar with the sequence con-	
	tent. Changes in perplexity can reflect the extent to	
	which a model is familiar with a particular linguis-	
	tic style. Using the Base SFT model as a fixed eval-	
	uator, we compute changes in the PPL of outputs	
	produced under different training methods; these	
	changes then reflect stylistic shifts during training.	
	During judging stage training of JudgeRLVR, we	
	sample model outputs at each step (10 samples per	
	step) for Vanilla RLVR and JudgeRLVR, and com-	
	pute the average PPL of the Base SFT model on	
	these outputs.	

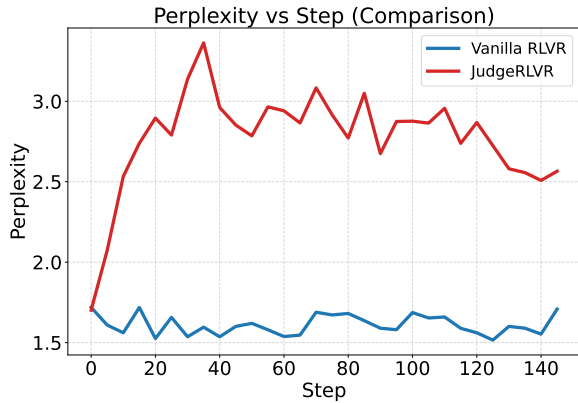


Figure 3: Base SFT perplexity (PPL) evaluated on sampled outputs along training steps. Vanilla RLVR stays stylistically close to Base SFT (flat PPL), while JudgeRLVR (judging stage) exhibits increasing PPL, indicating style transfer.

Transition-word statistics. Behaviors such as explicit backtracking and reflection in model outputs are typically manifested through the use of transition-words. Therefore, by counting the number and frequency of such discourse markers in outputs over the course of training, these statistics can indicate, to some extent, that backtracking is decreasing. During generating stage training of JudgeRLVR, we sample 1000 outputs per step and count occurrences and frequencies of contrast/backtracking markers:

but, however, though, although, yet,
 still, nevertheless, nonetheless,
 instead, conversely, whereas, while,
 actually, wait.

The above experiments are used to demonstrate shifts in style and backtracking. Combined with improvements in model performance, they can provide a comprehensive indication that these shifts are moving in a favorable direction.

5 Results

5.1 Main Results

Table 1 compares Base SFT, Vanilla RLVR that optimizes only final-answer correctness, and our proposed JudgeRLVR (sequential two-stage). The overall conclusion is that, compared with Vanilla RLVR, JudgeRLVR achieves higher or comparable accuracy while producing substantially shorter generations on most benchmarks, demonstrating a better quality–efficiency trade-off.

In-domain math: improving accuracy while markedly reducing redundant reasoning.

On AIME24/25, HMMT_feb_2025, and BeyondAIME, JudgeRLVR consistently improves accuracy over Vanilla RLVR, while significantly reducing generation length. This matches the motivation of our divided stages: judging stage, via strict balanced sampling (equal numbers of correct and incorrect traces per problem), strengthens the model’s ability to recognize error patterns and unproductive reasoning; then in generating stage, generating optimization can more reliably follow higher-confidence, lower-branching solution paths, reducing trial-and-error backtracking and verbose exploration. The larger gains on HMMT and BeyondAIME further suggest that the advantage of “judge first, then generate” is more pronounced when problems require longer chains of reasoning and stronger strategy selection.

On MATH500, JudgeRLVR is slightly worse in accuracy than Vanilla RLVR, but reduces generation length drastically. This indicates that for relatively saturated datasets, there is limited headroom for further accuracy gains, while JudgeRLVR can still substantially reduce reasoning cost when accuracy is largely maintained.

Out-of-domain: both generalization and efficiency gains, with task-dependent trade-offs.

On GPQA Diamond, LiveCodeBenchv6, and MMLU-Redux, JudgeRLVR improves accuracy while also reducing length. This suggests that the “error-sensitive / avoid-detours” preference learned in judging stage is not restricted to math, but can transfer to science reasoning, coding, and broad knowledge tasks, reducing ineffective exploration and improving overall decision quality.

In contrast, on IFEval and ZebraLogic we observe “higher accuracy but longer outputs”. This suggests that for tasks emphasizing format compliance, constraint satisfaction, or discrete rule verification, the model may need to generate more explicit checks to ensure correctness, trading off some length advantage for a higher pass rates. This does not contradict our goal: JudgeRLVR is not designed to mechanically minimize length, but to reduce *unproductive* reasoning; when a task intrinsically requires more explicit structure, a moderate increase in length may be a necessary cost.

5.2 Ablation Results

Table 2 reports two ablations—Judge Only and Mixed Strategy—to test: (i) whether judging stage alone is sufficient; and (ii) whether sequential stag-

Table 1: Main results comparing Base SFT, Vanilla RLVR, and JudgeRLVR (Sequential). We report accuracy (Acc) and average generation length in tokens (Len). For JudgeRLVR, Δ is absolute change in accuracy (pp:percentage point) vs Vanilla RLVR (higher is better); for Len, Δ is relative change vs Vanilla RLVR (shorter is better).

Benchmark	Base SFT		Vanilla RLVR		JudgeRLVR (Seq)			
	Acc \uparrow	Len \downarrow	Acc \uparrow	Len \downarrow	Acc \uparrow	Δ (pp)	Len \downarrow	Δ (%)
<i>In-Domain (Math)</i>								
AIME24	85.7	20.8k	86.3	21.8k	89.0	+2.7	12.9k	-41
AIME25	74.1	26.8k	77.0	27.2k	78.7	+1.7	17.2k	-37
MATH500	97.4	6.7k	98.0	6.8k	97.2	-0.8	2.0k	-71
HMMT_feb_2025	57.1	24.9k	60.8	26.9k	70.0	+9.2	20.4k	-24
BeyondAIME	55.5	33.8k	58.3	34.3k	63.9	+5.6	21.3k	-38
Overall Avg	74.0	22.6k	76.1	23.4k	79.8	+3.7	14.8k	-42
<i>Out-of-Domain</i>								
GPQA Diamond	63.0	12.0k	61.2	12.1k	66.4	+5.2	11.2k	-7.5
IFEval	80.6	3.2k	79.9	4.1k	86.4	+6.5	4.6k	+12
LiveCodeBenchv6	57.4	20.6k	58.2	22.8k	63.9	+5.7	18.7k	-18
MMLU Redux	86.8	2.1k	85.4	2.1k	87.0	+1.6	1.5k	-28
ZebraLogic	84.6	7.7k	83.6	7.9k	86.8	+3.2	9.6k	+22
Overall Avg	74.5	9.1k	73.7	9.8k	78.1	+4.5	9.1k	-3.9

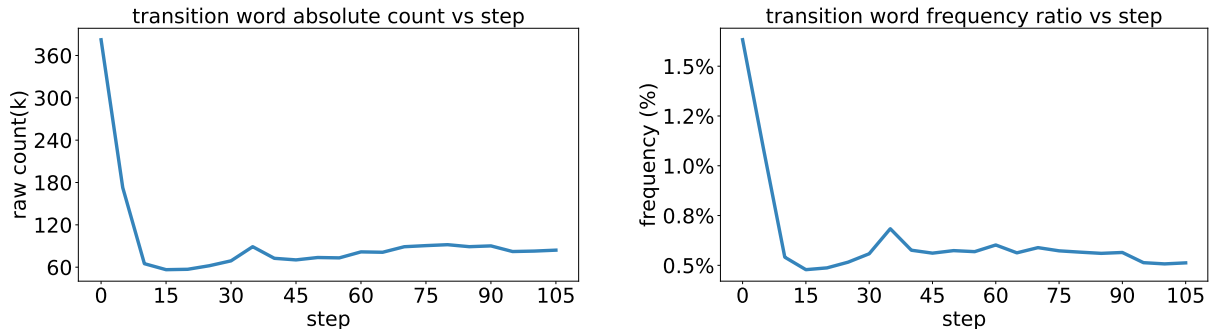


Figure 4: Counts(left) and Frequencies(right) of transition/backtracking markers (e.g., *but*, *however*, *wait*) in sampled outputs across JudgeRLVR generating stage training steps.

482 ing is essential.

483 **Judge Only: better “judging” does not reliably**
 484 **translate into better generation.** Judge Only
 485 shows clear degradations on multiple benchmarks,
 486 especially on in-domain math: compared with
 487 JudgeRLVR, accuracy drops on math benchmarks.
 488 More importantly, generation length increases sub-
 489 stantially. This indicates that the discriminative
 490 capability learned in judging stage does not auto-
 491 matically induce a more concise generation policy;
 492 it can even make the model more prone to “writing
 493 out the checking process” or reasoning more cau-
 494 tiously and verbosely, which lengthens outputs and
 495 hurts final accuracy. In other words, error-aware

496 discrimination is a necessary prerequisite, but not a
 497 substitute for generating optimization; generating
 498 stage is still required to convert error sensitivity
 499 into better path selection and more information-
 500 dense expression during generation.

501 We also observe that Judge Only can achieve
 502 comparable or slightly better accuracy on some
 503 out-of-domain tasks while producing shorter out-
 504 puts. This suggests that some “judging” prefer-
 505 ences learned in judging stage may directly benefit
 506 certain tasks; however, this benefit is unstable and
 507 cannot replace the full two-stage pipeline for our
 508 core objective.

Table 2: Main results comparing JudgeRLVR (Sequential), Judge Only, and Mixed Strategy. We report accuracy (Acc) and average generation length in tokens (Len). For Judge Only and Mixed Strategy, Δ is absolute change in accuracy (pp:percentage point) vs JudgeRLVR (higher is better); for Len, Δ is relative change vs JudgeRLVR (shorter is better).

Benchmark	JudgeRLVR (Seq)		Judge Only				Mixed Strategy			
	Acc \uparrow	Len \downarrow	Acc \uparrow	Δ (pp)	Len \downarrow	Δ (%)	Acc \uparrow	Δ (pp)	Len \downarrow	Δ (%)
<i>In-Domain (Math)</i>										
AIME24	89.0	12.9k	87.0	-2.0	22.4k	+74	88.3	-0.7	12.8k	-0.8
AIME25	78.7	17.2k	76.3	-2.4	28.5k	+66	79.7	+1.0	18.2k	+6.0
MATH500	97.2	2.0k	97.2	+0.0	7.4k	+266	98.0	+0.8	2.2k	+7.9
HMMT_feb_2025	70.0	20.4k	61.5	-8.5	17.9k	-12	69.0	-1.0	25.9k	+27
BeyondAIME	63.9	21.3k	57.3	-6.6	35.8k	+68	63.4	-0.5	22.9k	+7.3
Overall Avg	79.8	14.8k	75.9	-3.9	22.4k	+92	79.7	-0.1	16.4k	+9.5
<i>Out-of-Domain</i>										
GPQA Diamond	66.4	11.2k	66.9	+0.5	7.9k	-30	63.7	-2.7	16.1k	+44
IFEval	86.4	4.6k	85.9	-0.5	4.2k	-9.0	80.4	-6.0	7.9k	+72
LiveCodeBenchv6	63.9	18.7k	61.4	-2.5	21.3k	+14	56.4	-7.5	26.7k	+43
MMLU Redux	87.0	1.5k	87.4	+0.4	1.0k	-31	87.1	+0.1	3.3k	+119
ZebraLogic	86.8	9.6k	87.3	+0.5	8.0k	-17	87.4	+0.6	7.9k	-19
Overall Avg	78.1	9.1k	77.8	-0.3	8.5k	-15	75.0	-3.1	12.4k	+52

Mixed Strategy: interleaving judge and generate updates weakens the benefit of sequential stages.

Mixed Strategy can be close to JudgeRLVR on some tasks, but is overall more unstable: it slightly exceeds JudgeRLVR on AIME25 and MATH500, yet drops substantially on tasks such as IFEval and LiveCodeBenchv6, often with significantly longer outputs. This aligns with our hypothesis: when “judge-style” and “generation-style” behaviors are optimized in a 1:1 interleaved manner, the model must satisfy two different objectives and reward structures in the same phase, which can lead to interference in policy-gradient updates and prevent the emergence of a clean, stable internal decision process of “judge first, then generate”. The result is occasional local accuracy gains, but weaker overall efficiency and generalization.

5.3 Mechanism Verification

Style transfer measured by PPL under Base SFT.

Figure 3 reports the Base SFT model’s average perplexity on sampled outputs across training steps. For Vanilla RLVR, the average PPL remains largely unchanged, suggesting the output distribution stays stylistically close to Base SFT. In contrast, during JudgeRLVR judging stage, Base SFT PPL increases noticeably, indicating a systematic style shift induced by judge training.

Reduced explicit backtracking measured by transition words.

Figure 4 show that during

JudgeRLVR generating stage, both the absolute count and frequency of transition/backtracking words decrease substantially over training, consistent with reduced explicit backtracking and a shift of verification/correction into internal decision-making.

6 Conclusion

In this paper, we proposed JudgeRLVR, a two-stage training paradigm that acts as a precursor to efficient reasoning. By explicitly training the model to discriminate solution validity before generating, we achieved a superior quality–efficiency trade-off. Our experiments on Qwen3-30B-A3B show that JudgeRLVR improves in-domain math accuracy while reducing generation length by approximately 42%, with notable generalization to out-of-domain tasks.

Many assume that judging an answer is the easy part, because verification seems to come with a checklist that generation does not. Our study points to a sharper truth: on identical math problems with the same model, judging yields greater improvement than solving alone, so “is this correct?” is a problem in its own right. Training the judge makes the model internalize what a good line of thought looks like, and later generation can cross its own limits with fewer false starts. In this sense, to generate is simply to judge so well that the correct answer is the only thing left to say.

567	Limitations		
568	Despite achieving a favorable quality–efficiency		
569	trade-off on RLVR and showing strong generaliza-		
570	tion, our work has several limitations that merit		
571	further investigation:		
572	Limited training domains Our experiments only		
573	train on mathematics-domain data. Whether the		
574	same training strategy and observed gains transfer		
575	to other domains (e.g., coding, scientific QA, open-		
576	domain reasoning) remains to be validated.		
577	Coarse-grained judge signals We employ an		
578	outcome-level judge on the final solution response		
579	with a binary correct/incorrect decision. This super-		
580	vision may be insufficient to capture fine-grained		
581	reasoning errors (e.g., identifying the exact incor-		
582	rect step), and more granular judge feedback could		
583	further improve learning and diagnostic ability.		
584	References		
585	Mislav Balunović, Jasper Dekoninck, Ivo Petrov, Nikola		
586	Jovanović, and Martin Vechev. 2025. Matharena: Evaluating llms on uncontaminated math competitions.		
587			
588			
589	[ByteDance-Seed]. 2025. Beyondaime: Advancing math reasoning evaluation beyond high school olympiads. https://huggingface.co/datasets/ByteDance-Seed/BeyondAIME.		
590			
591			
592			
593			
594			
595	Hongyi James Cai and et al. 2025. How much backtracking is enough? exploring the interplay of sft and rl in enhancing llm reasoning.		
596			
597			
598	Micheline T. H. Chi, Paul J. Feltovich, and Robert		
599	Glaser. 1981. Categorization and representation of		
600	physics problems by experts and novices. <i>Cognitive</i>		
601	<i>Science</i> , 5(2):121–152.		
602	AI DataLearner. 2025. Reinforcement learning with verifiable rewards (rlvr): Why the focus of llm training is shifting in 2025. <i>DataLearner Official Blog</i> .		
603			
604			
605	Aryo P. Gema, Joshua O. J. Leang, Giwon Hong,		
606	Alessio Devoto, Alberto C. M. Mancino, Rohit Sax-		
607	ena, Xuanli He, Yu Zhao, Xiaotang Du, Moham-		
608	mad R. G. Madani, Claire Barale, Robert McHardy,		
609	Joshua Harris, Jean Kaddour, Emile Van Krieken,		
610	and Pasquale Minervini. 2025. Are we done with mmlu? In <i>Proceedings of NAACL 2025</i> .		
611			
612	Jian Guo, Jie Tang, and Xiaolin Li. 2025a. Deepseek-r1: Scaling rlvr for mathematical reasoning.		
613			
614	Yuchen Guo, Xin Wang, Yifan Zhang, and Pengfei Liu.		
615	2025b. Discriminative finetuning of generative large language models without reward models and preference data.	616	617
	Xu Han, Yan Wang, Zonghan Li, and Chengqi Zhang.	618	619
	2024. Uncertainty-based verifiers for llm reasoning outputs.	620	
	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul	621	622
	Arora, Steven Basart, Eric Tang, Dawn Song, and	623	624
	Jacob Steinhardt. 2021. Measuring mathematical problem solving with the MATH dataset. <i>Advances in Neural Information Processing Systems (NeurIPS)</i> .	625	
	Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia	626	627
	Yan, Tianjun Zhang, Sida Wang, Armando Solar-	628	629
	Lezama, Koushik Sen, and Ion Stoica. 2024. Livecodebench: Holistic and contamination free evaluation of large language models for code.	630	
	Nathan Lambert, Noah A. Smith, and Hannaneh Ha-	631	632
	jishirzi. 2024. Scalable verifiable reward learning for large language models.	633	
	Bill Yuchen Lin, Ronan Le Bras, Kyle Richardson,	634	635
	Ashish Sabharwal, Radha Poovendran, Peter Clark,	636	637
	and Yejin Choi. 2025. Zebralogic: On the scaling limits of llms for logical reasoning.		
	Ilya Loshchilov and Frank Hutter. 2018. Decoupled Weight Decay Regularization. In <i>International Conference on Learning Representations</i> .	638	639
		640	
	Moonshot AI Team. 2025. Kimi k1.5: Scaling reinforcement learning with llms. <i>arXiv preprint arXiv:2501.12599</i> .	641	642
		643	
	David Rein, Betty Li Hou, Asa C. Stickland, Jackson	644	645
	Petty, Richard Y. Pang, Julien Dirani, Julian Michael,	646	647
	and Samuel R. Bowman. 2023. Gpqa: A graduate-level google-proof q&a benchmark.		
	Yuxin Shao, Tianyu Liu, Minghao Chen, and Shuohang	648	649
	Li. 2025. Spurious rewards: Rethinking training signals in rlvr.	650	
	Wei Shen, Jia Li, Hongyu Zhang, and Zhe Wang. 2025.	651	652
	Llm reasoning engine: Specialized training for en-	653	654
	hanced mathematical reasoning. In <i>Proceedings of the 4th International Workshop on Knowledge-Augmented Methods for Natural Language Processing (KnowledgeNLP’25)</i> , pages 118–128. Association for Computational Linguistics.	655	656
		657	
	Luis Shimabucoro, Rutuja Chavhan, Timothy	658	659
	Hospedales, and Henry Gouk. 2024. Evaluating the evaluators: Are validation methods for few-shot learning fit for purpose? <i>Transactions on Machine Learning Research</i> , pages 1–13.	660	661
		662	
	Zhiqing Sun, Yang Liu, Danqi Chen, and Wenhui Li.	663	664
	2024. Skywork-math: Scaling mathematical reasoning with corpus fine-tuning.	665	
	Xinyi Tang, Jiaqi Chen, Bing Zhang, and Zhiyuan Liu.	666	667
	2025. Unifiedreward-think: A comprehensive evaluator for reasoning quality.	668	

669 DeepSeek Team, Feng Liu, Hao Wang, and Bo Zhang.
670 2024. Deepseek-math: Enhancing mathematical rea-
671 soning via rlvr. In *Proceedings of the 38th Interna-*
672 *tional Conference on Machine Learning (ICML'24)*,
673 pages 2345–2356. PMLR.

674 Qwen Team. 2025. *Qwen3 technical report*. *Preprint*,
675 arXiv:2505.09388.

676 LLM-Core-Team Xiaomi. 2025. *Mimo: Unlocking the*
677 *reasoning potential of language model – from pre-*
678 *training to posttraining*. *Preprint*, arXiv:2505.07608.

679 Yilun Xu, Tuhin Chakraborty, Shubham Sharma, Lucas
680 Nunes, Emre Kıcıman, Shan Lu, and Rakesh Chan-
681 dra. 2025. *Direct reasoning optimization: Llms can*
682 *reward and refine their own reasoning for open-ended*
683 *tasks*.

684 Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan,
685 Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan,
686 Gaohong Liu, Lingjun Liu, and 1 others. 2025. *Dapo:*
687 *An open-source llm reinforcement learning system*
688 *at scale*. *arXiv preprint arXiv:2503.14476*.

689 Xiang Yue, Zheqi Chen, and Gao Huang. 2025. *Does*
690 *reinforcement learning with verifiable rewards elicit*
691 *novel reasoning abilities?*

692 Yifan Zhang and Team Math-AI. 2024. American invi-
693 tational mathematics examination (aime) 2024.

694 Yifan Zhang and Team Math-AI. 2025. American invi-
695 tational mathematics examination (aime) 2025.

696 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha
697 Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and
698 Le Hou. 2023. *Instruction-following evaluation for*
699 *large language models*.

700 Mingsheng Zhou, Jun Li, Xiaoyu Wang, and Lingling
701 Zhang. 2025. *Jiuzhang3.0: Efficient mathematical*
702 *reasoning via synthetic data pre-training*.

703 Xinyu Zhu, Mengzhou Xia, Zhepei Wei, Wei-Lin Chen,
704 Danqi Chen, and Yu Meng. 2025. *The surprising*
705 *effectiveness of negative reinforcement in llm reason-*
706 *ing*.

A Prompt Templates of Judge and Generate

Judging Prompt

Given the question:
{question}
And the proposed answer:
{solution response}

Analyze concisely whether the answer is correct. After your analysis, put only your final judgment in \boxed{} as 1 for correct and 0 for incorrect.

Generating and Rollout Prompt

{question}
Let's think step by step and output the final answer within \boxed{}

B Case Study

Full output of Figure 1

Problem Statement

Convert the point $(0, 3)$ in rectangular coordinates to polar coordinates. Enter your answer in the form (r, θ) , where $r > 0$ and $0 \leq \theta < 2\pi$.

Vanilla RLVR

<think>

I need to convert the point $(0,3)$ from rectangular coordinates to polar coordinates. The rectangular coordinates are given as $(x,y) = (0,3)$, and I need to find the polar coordinates (r, θ) , where r is the distance from the origin, and θ is the angle from the positive x-axis.

First, recall the formulas for conversion. The radius r is given by $r = \sqrt{x^2 + y^2}$. Plugging in the values, x is 0 and y is 3, so $r = \sqrt{0^2 + 3^2} = \sqrt{0 + 9} = \sqrt{9} = 3$. So r is 3.

Now, for the angle θ , it's given by $\theta = \tan^{-1}(y/x)$. But here, x is 0, and y is 3. So if I plug in, $\tan^{-1}(3/0)$, which is undefined because division by zero. That makes sense because when x is 0, the point is on the y-axis.

I know that when a point is on the positive y-axis, the angle θ should be $\pi/2$ radians, or 90 degrees. Since the problem asks for θ in radians between 0 and 2π , I should use radians.

So, $\theta = \pi/2$.

Let me confirm. If r is 3, and θ is $\pi/2$, then the rectangular coordinates should be $x = r \cos \theta = 3 * \cos(\pi/2) = 3 * 0 = 0$, and $y = r \sin \theta = 3 * \sin(\pi/2) = 3 * 1 = 3$. Perfect, that matches.

The problem says $r > 0$, which is satisfied since $r=3>0$, and θ between 0 and 2π , $\pi/2$ is within that range.

I recall that sometimes people might think of θ as $-\pi/2$ or $3\pi/2$, but that would give a negative y-value, which is not the case here. Since y is positive, it has to be in the first or second quadrant, but since x is 0, it's specifically on the positive y-axis, so $\theta = \pi/2$.

If θ were $3\pi/2$, that would be the negative y-axis, where y is negative, but here y is positive, so it's not that.

So, I think $(3, \pi/2)$ is the correct polar coordinate.

The answer should be in the form (r, θ) , so $(3, \pi/2)$.

762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812

I should write it as a pair. And the problem says to box the final answer.

So, I think that's it.

Just to be thorough, is there any other representation? For example, could r be negative? But the problem specifies $r > 0$, so no. And θ is usually taken in $[0, 2\pi)$, so $\pi/2$ is fine.

Sometimes people use degrees, but here it's specified to use radians, and the range is given, so radians are implied.

The problem doesn't specify units, but in mathematics, especially with π , it's standard to use radians.

So, I'm confident.

</think>

First, I need to convert the point $(0, 3)$ from rectangular coordinates to polar coordinates (r, θ) with $r > 0$ and $0 \leq \theta < 2\pi$, use the conversion formulas.

The radius r is given by:

$$r = \sqrt{x^2 + y^2} = \sqrt{0^2 + 3^2} = \sqrt{9} = 3.$$

The angle θ is found using $\theta = \tan^{-1}(y/x)$. However, since $x = 0$ and $y = 3 > 0$, the point lies on the positive y -axis. Thus, $\theta = \pi/2$.

Verify the solution:

- $x = r \cos \theta = 3 \cos(\pi/2) = 3 \cdot 0 = 0$,
- $y = r \sin \theta = 3 \sin(\pi/2) = 3 \cdot 1 = 3$,

which matches the original point.

The conditions are satisfied: $r = 3 > 0$ and $\theta = \pi/2$ is in $[0, 2\pi)$.

Thus, the polar coordinates are $(3, \pi/2)$.

The final answer is:

$$\left(3, \frac{\pi}{2}\right)$$

JudgeRLVR

<think>

We need to convert $(0,3)$ to polar coordinates. $r = \sqrt{x^2 + y^2} = \sqrt{0+9}=3$. Theta is angle from positive x -axis to point $(0,3)$. That's 90 degrees = $\pi/2$ rad. Since $x=0, y>0$, theta = $\pi/2$. So answer $(3, \pi/2)$. Ensure within 0 to 2π . So final answer

$$\left(3, \frac{\pi}{2}\right)$$

</think>

First, I need to calculate r .

The point $(0, 3)$ has

$$r = \sqrt{0^2 + 3^2} = 3$$

and lies on the positive y -axis, so the angle from the positive x -axis is

$$\theta = \frac{\pi}{2}.$$

Thus the polar coordinates are $(3, \pi/2)$.

The final answer is:

$$\left(3, \frac{\pi}{2}\right)$$

Table 3: Evaluation Benchmarks

Benchmark	Metric
AIME24	Acc avg@10
AIME25	Acc avg@10
MATH500	Acc
HMMT_feb_2025	Acc avg@10
BeyondAIME	Acc avg@8
GPQA Diamond	Exact Match avg@5
IFEval	Prompt Strict
LiveCodeBenchv6	Pass@1
MMLU-Redux	Exact Match
ZebraLogic	Acc

C Evaluate Details

Table 3 shows the metric setup of evaluation benchmarks.

813
814
815