

Semantic-Space Exploration and Exploitation in RLVR for LLM Reasoning

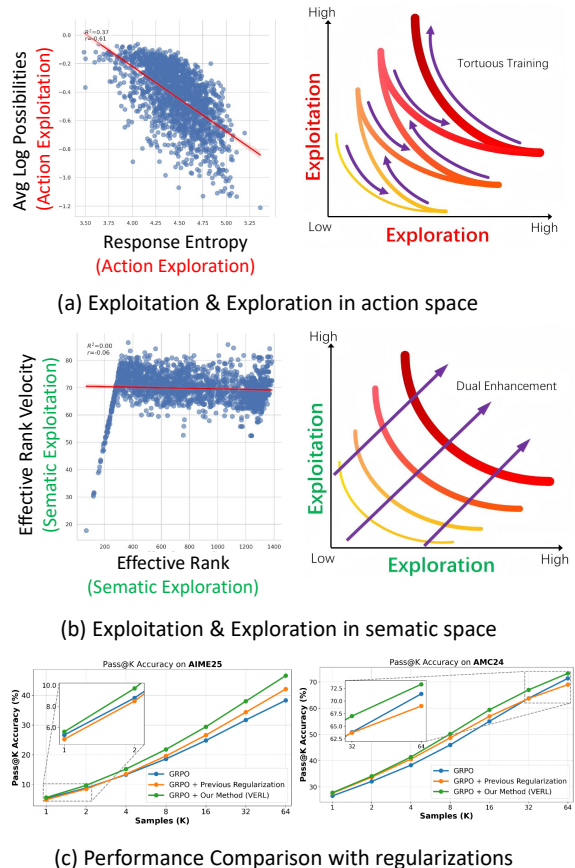
Anonymous ACL submission

Abstract

Reinforcement Learning with Verifiable Rewards (RLVR) for LLM reasoning is often framed as balancing exploration and exploitation in action space, typically operationalized with token-level proxies (e.g., output entropy or confidence). We argue that this apparent trade-off is largely a measurement artifact: token-level statistics reflect next-token uncertainty rather than how reasoning progresses over multi-token semantic structures. We therefore study exploration and exploitation in the hidden-state space of response trajectories. We use Effective Rank (ER) to quantify representational exploration and introduce its temporal derivatives, Effective Rank Velocity (ERV) and Effective Rank Acceleration (ERA), to characterize exploitative refinement dynamics. Empirically and theoretically, ER and ERV exhibit near-zero correlation in semantic space, suggesting the two capacities can be improved simultaneously. Motivated by this, we propose Velocity-Exploiting Rank Learning (VERL), which shapes the RL advantage with an auxiliary signal derived from ER/ERV and uses the more stable ERA as a meta-control variable to adaptively balance the incentives. Across multiple base models, RL algorithms, and reasoning benchmarks, VERL yields consistent improvements, including large gains on challenging tasks (e.g., 21.4% in Gaokao 2024).

1 Introduction

Recent advancements in Reinforcement Learning with Verifiable Rewards (RLVR) have substantially improved the reasoning abilities of Large Language Models (LLMs). A dominant narrative in recent work (Chen et al., 2025; Yue et al., 2025; Deng et al., 2025a; Agarwal et al., 2025) interprets this progress through the lens of balancing *exploration* (seeking diverse reasoning paths) and *exploitation* (refining the most promising known strategies). In practice, however, these notions are almost exclusively operationalized in the *token-level*



(a) Exploitation & Exploration in action space
(b) Exploitation & Exploration in semantic space
(c) Performance Comparison with regularizations

Figure 1: Comparative analysis with the responses of DeepSeek-R1-Distill-Qwen-7B in simpleRL test dataset (Zeng et al., 2025). (a) Traditional metrics for exploitation & exploration constrained by negative coupling, leading to meandering progress for both capabilities. (b) Our hidden-state metrics are structurally decoupled, capturing complementary aspects of semantic exploration and exploitation. (c) Training regularization with our metrics demonstrates stronger performance in both exploitation (small K) and exploration (large K).

action space: exploration is associated with high-entropy next-token distributions, while exploitation is linked to high-confidence ones. These token-level proxies are convenient to compute, yet they primarily measure uncertainty about the *next token*. As a result, the widely discussed “exploration–exploitation trade-off” (Fig. 1a) in reasoning may

reflect the limitations of these proxies, instead of an inherent property of reasoning behavior.

This mismatch becomes more evident when distinguishing the *action space* from the *semantic space* in LLM reasoning. Reasoning progress is better viewed as transitions between *semantic states*, traversing concepts, subgoals, and partial derivations that span many tokens (Wei et al., 2022; Yao et al., 2023) while token sampling occurs in the action space. Token-level proxies (e.g., next-token entropy or confidence) can therefore misattribute local sampling uncertainty to semantic progress, leading to two failure modes: higher entropy may only amplify lexical variability without reaching new semantic directions, whereas higher confidence may collapse the trajectory too early into a familiar semantic basin (Fu et al., 2025; Qiao et al., 2025; Agarwal et al., 2025).

As prior work has shown that pre-trained LLM hidden states naturally encode rich semantic structure (Huang et al., 2025; Jing et al., 2025) and reasoning dynamics (Cheng et al., 2025a; Deng et al., 2025b), they can provide a continuous representation of the model’s evolving semantic state. This suggests that exploration and exploitation should be defined *over trajectories in semantic space* (how broadly the representation spans semantic directions, and how it refines within a region), rather than over token distributions in action space. We therefore ask: *Is the exploration–exploitation trade-off intrinsic to reasoning, or largely an artifact of measuring it in the token-level action space?*

To answer this, we step out of the action space and study exploration and exploitation directly in a *semantic space* induced by transformer hidden-state representations. In this semantic space, we quantify *representational exploration* by Effective Rank (ER), which measures how broadly a hidden-state trajectory spans semantic directions (Valeriani et al., 2023; Matthews et al., 2024; Jing et al., 2025). Following this intuition, we quantify *representational exploitation* by the temporal evolution of this span: Effective Rank Velocity (ERV) captures how semantic complexity is refined along the trajectory, and Effective Rank Acceleration (ERA) captures whether this refinement is sustaining, accelerating, or saturating. Equipped with these semantic-trajectory measures, we uncover a striking empirical result: ER and ERV exhibit near-zero correlation (Fig. 1b). This provides evidence that the commonly perceived exploration–exploitation trade-off

is not intrinsic to LLM reasoning, but largely an artifact of measuring both notions in the token-level action space. And the two capacities can be *decoupled and improved simultaneously* in semantic space (Fig. 1c).

Building on this insight, we propose Velocity-Exploiting Rank Learning (VERL), a plug-and-play advantage-shaping method that optimizes exploration and exploitation *in semantic space*. VERL augments the RL advantage with auxiliary incentives derived from ER (encouraging broader semantic directions to avoid premature collapse) and ERV (reinforcing productive semantic refinement). Crucially, VERL uses ERA as a meta-control signal to adapt the strength of these incentives over training, as ERA tracks the trend of semantic refinement (i.e., whether exploitation is accelerating or saturating) and, moreover, is less sensitive to short-term fluctuations with theoretical stability guarantees (Sec. 2.3). Across diverse base models, RL algorithms, and reasoning benchmarks, VERL consistently improves performance, including up to **21.4%** absolute accuracy on the challenging Gaokao 2024 benchmark.

Contributions. (i) We move exploration/exploitation from the token-level action space to a hidden-state semantic space, introducing ER, ERV, and ERA to measure semantic trajectories. (ii) We find that ER and ERV are nearly uncorrelated, indicating that the commonly reported trade-off is largely influenced by token-level measurement choices. (iii) We present VERL, an advantage-shaping method that improves both capacities, and show consistent gains across models, RL algorithms, and benchmarks.

2 Preliminaries

2.1 Problem Formulation and Notations

We adopt an RLVR perspective for LLM reasoning. Given a prompt $x \sim \mathcal{P}_x$, an LLM policy $\pi_\theta(\cdot | x)$ generates a trajectory $y_{0:T} = (y_0, \dots, y_T)$ with $T \leq L_{\max}$. A scalar reward $r(x, y)$ evaluates the completed trajectory. We learn parameters θ by maximizing the expected reward:

$$\phi = \operatorname{argmax}_{\theta} \mathbb{E}_{x \sim \mathcal{P}_x} \mathbb{E}_{y \sim \pi_\theta(\cdot | x)} [r(x, y)]. \quad (1)$$

A common view is that optimizing this objective entails an exploration–exploitation trade-off: the policy must explore diverse reasoning paths to discover high-reward solutions while exploiting reliable strategies that consistently succeed.

2.2 Reinforcement Learning Baselines

Proximal Policy Optimization (PPO) from Schulman et al. (2017) is a standard RL algorithm that seeks to maximize a clipped surrogate objective function. This objective prevents excessively large changes that would destabilize training, which is defined as:

$$\mathcal{L}_{\text{PPO}}(\theta) = \mathbb{E}_{x,y} \left[\sum_{t=1}^{|y|} \min(\rho_t(\theta)A_t, \bar{\rho}_t(\theta)A_t) \right],$$

$$\bar{\rho}_t(\theta) := \text{clip}_{1-\epsilon_{\text{low}}}^{1+\epsilon_{\text{high}}}(\rho_t(\theta)), \quad (2)$$

where $\rho_t(\theta) := \frac{\pi_{\theta}(y_t|x, y_{<t})}{\pi_{\theta_{\text{old}}}(y_t|x, y_{<t})}$ is the probability ratio between the current and old policies, and A_t is the estimated advantage, often calculated using Generalized Advantage Estimation (GAE) from Schulman et al. (2016), with clipping (hyperparameter ϵ) to mitigate excessive deviation.

Group Relative Policy Optimization (GRPO) from Shao et al. (2024) computes a baseline directly from the rewards of multiple trajectories. For a given prompt, it samples a group of G responses, obtains their corresponding rewards $\{r_1, \dots, r_G\}$, and normalizes these rewards to compute the advantage for each response:

$$A_{i,t} := (r_i - \text{mean}(\{r_j\}_{j=1}^G)) / \text{std}(\{r_j\}_{j=1}^G). \quad (3)$$

GRPO would assign a single rule-based reward to the entire output sequence, and the resulting group-relative advantage is uniformly propagated to all tokens, then updated as in Eq. 2.

2.3 Hidden State Representations

Response hidden states. For an autoregressive response of length T , let $z_t \in \mathbb{R}^D$ denote the final-layer hidden state at token t . Stacking them yields $\mathbf{Z} = [z_1; \dots; z_T] \in \mathbb{R}^{T \times D}$, which we view as a semantic trajectory. **Dataset hidden states.** Given N prompts, we summarize the i -th response by $\bar{z}_i = \frac{1}{T_i} \sum_{t=1}^{T_i} z_{i,t}$ (Skean et al., 2025), and stack them into $\bar{\mathbf{Z}} = [\bar{z}_1; \dots; \bar{z}_N] \in \mathbb{R}^{N \times D}$ to represent the dataset-level semantic distribution.

3 A Hidden-State Perspective on Representational Dynamics

3.1 Semantic Exploration: Effective Rank (ER)

Following Roy and Vetterli (2007), we define the Effective Rank of a response hidden-state matrix $\mathbf{Z} \in \mathbb{R}^{T \times D}$ using its (non-padding) singular values

$\{\sigma_j\}$. Let $p_j = \sigma_j / \sum_k \sigma_k$. Then

$$\text{ER}(\mathbf{Z}) := \text{erank}(\mathbf{Z}) = \exp \left[- \sum_j p_j \log p_j \right]. \quad (4)$$

We interpret ER as *representational breadth* in the hidden-state semantic space: larger ER indicates that the trajectory spans a broader set of semantic directions, while smaller ER suggests concentration in fewer directions.

Theorem 3.1. *Suppose we have a matrix of embeddings $\mathbf{Z} \in \mathbb{R}^{T \times D}$. Then the ER of \mathbf{Z} is a lower bound of conventional rank of \mathbf{Z} :*

$$1 \leq \text{erank}(\mathbf{Z}) \leq \text{rank}(\mathbf{Z}) \leq \min\{T, D\}. \quad (5)$$

Remark 3.2. While rank is discrete and blind to spectral concentration, ER is a continuous effective-dimensionality measure that grows as the singular spectrum flattens. Details are in App. E.2.

3.2 Semantic Exploitation: Effective Rank Velocity (ERV) and Acceleration (ERA)

Effective Rank (ER) provides a *static* snapshot of how broadly a hidden-state trajectory spans semantic directions. To capture *how this semantic breadth evolves during reasoning*, namely whether the trajectory continues to expand into new directions, steadily refines within a region, or gradually saturates, we introduce two temporal-difference metrics: Effective Rank Velocity (ERV) and Effective Rank Acceleration (ERA). Importantly, these dynamics are defined *in semantic space* through hidden-state representations rather than in the token-level action space. This choice aligns with our view that exploitation reflects *productive semantic refinement*, not merely low-entropy sampling.

Setup. Let M be any trajectory metric computed on hidden states (we will instantiate $M = \text{ER}$). Given a response of length T , let m_t denote the value of M computed on the prefix hidden-state matrix up to position t . To reduce sensitivity to short-range token fluctuations and control computation, we evaluate the metric every s tokens. Define the evaluation indices $\mathcal{T} = \{s, 2s, \dots, Ks\}$, where $K = \lfloor (T-1)/s \rfloor$.

Definition 3.3 (ERV: First-order temporal difference). For each evaluation step $j \in \{2, \dots, K\}$, we define an *instantaneous improvement* term

$$\delta_{j \cdot s} := m_{j \cdot s} - \frac{1}{j-1} \sum_{k=1}^{j-1} m_{k \cdot s}. \quad (6)$$

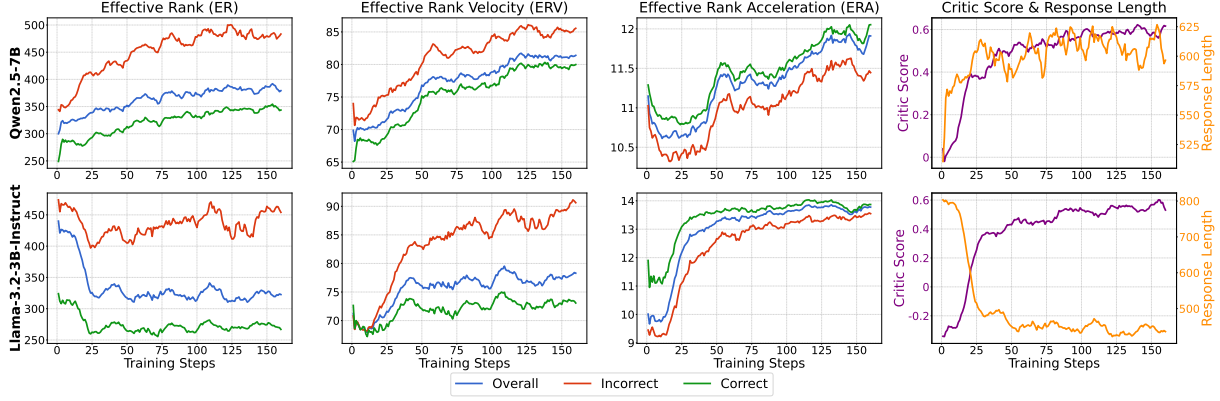


Figure 2: Response-level metrics during GRPO post-training, smoothed with a 10-step rolling window. Metrics are shown for the Overall batch, as well as for subsets of Correct and Incorrect samples. The rightmost column displays the average Critic Score (reward) and Response Length per batch.

This contrasts the current semantic breadth with the historical average over earlier prefixes. We then aggregate these terms to obtain the first-order temporal difference:

$$\Delta_M^{(1)} := \frac{1}{K-1} \sum_{j=2}^K \delta_{j \cdot s}. \quad (7)$$

ERV tracks net semantic progress against running history, aligning with exploitation as greedy trajectory refinement under the current policy guidance (App. E.3). A large positive $\delta_{j \cdot s}$ indicates that the trajectory at step $j \cdot s$ expands semantic breadth beyond its previous trend; a small (or negative) value suggests that semantic expansion is no longer outpacing the past. We instantiate our velocity metric as $\text{ERV} := \Delta_{\text{ER}}^{(1)}$, and interpret it as a measure of *exploitation in semantic space*: sustained positive ERV reflects consistent representational refinement along the current reasoning path, whereas small ERV indicates diminishing returns and potential stagnation.

Writing the consecutive-step increments as $\Delta m_{r \cdot s} := m_{r \cdot s} - m_{(r-1) \cdot s}$, we can rearrange $\delta_{j \cdot s}$ as

$$\delta_{j \cdot s} = \frac{1}{j-1} \sum_{r=2}^j (r-1) \Delta m_{r \cdot s}, \quad j \geq 2, \quad (8)$$

which shows that $\delta_{j \cdot s}$ is a time-weighted average of local increments, emphasizing recent changes.

Definition 3.4 (ERA: Second-order temporal difference). To capture whether the refinement dynamics are *accelerating* or *saturating*, we define the second-order temporal difference as

$$\Delta_M^{(2)} := \frac{1}{K-2} \sum_{j=3}^K [\delta_{j \cdot s} - \delta_{(j-1) \cdot s}]. \quad (9)$$

Interpretation. We instantiate $\text{ERA} := \Delta_{\text{ER}}^{(2)}$. Positive ERA indicates that the semantic refinement reflected by $\delta_{j \cdot s}$ is strengthening over time (accelerating refinement), while negative ERA suggests that refinement is weakening, consistent with stabilization or saturation in a narrow semantic region. In later sections, we leverage this stability-oriented signal to modulate training incentives in a principled and adaptive manner.

3.3 Scaling Properties of Semantic-Space Dynamics

In Sec. 3.1 and Sec. 3.2, we defined ER, ERV, and ERA on the hidden-state trajectory of a single response, i.e., in the model’s semantic space. We now study how these representational dynamics scale at two levels: (i) **dataset-level** scaling as the number of responses N grows, and (ii) **trajectory-level** scaling as the reasoning length T increases. The key observation is that both settings can be analyzed through the same lens: the effective number of approximately orthogonal semantic directions, denoted by k . The following proposition provides a unified model.

Proposition 3.5. *Assume a hidden-state matrix contains k approximately orthogonal row vectors. Then ER and its first-order temporal difference ERV scale linearly with k : $\text{ER} = \mathcal{O}(k)$ and $\Delta_{\text{ER}}^{(1)} = \mathcal{O}(k)$. In contrast, the second-order difference ERA does not scale with k : $\Delta_{\text{ER}}^{(2)} = \mathcal{O}(1)$.*

Remark 3.6. Prop. 3.5 admits two readings depending on what the rows of \mathbf{Z} represent. **Dataset level:** with $\mathbf{Z} = \bar{\mathbf{Z}} \in \mathbb{R}^{N \times D}$ (rows are response embeddings), k measures the number of distinct semantic modes covered by the dataset, so ER/ERV grow with dataset diversity while ERA stays $\mathcal{O}(1)$.

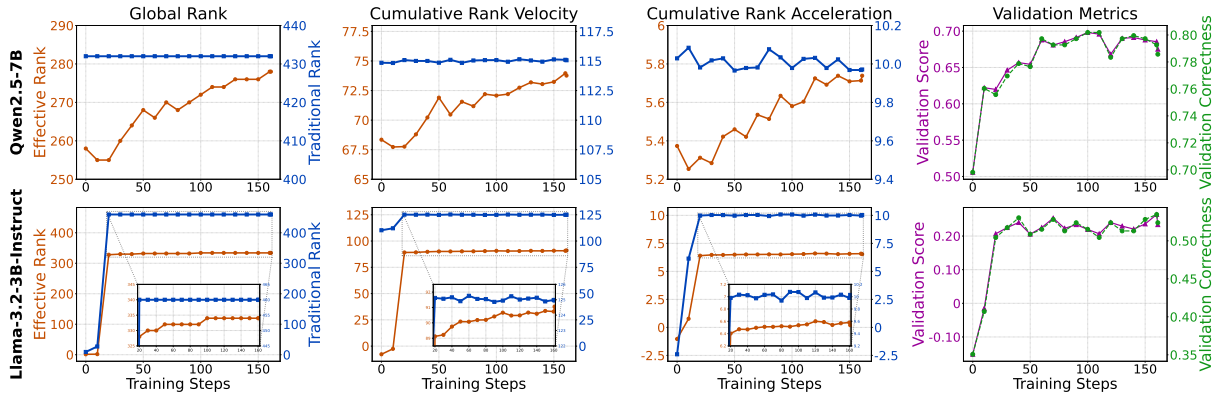


Figure 3: Visualization of dataset-level metrics during GRPO post-training. The figure compares **Traditional** metrics with our **proposed** metrics. Also shown are the **Validation Score** and sample **Correctness**, both averaged over the validation dataset.

Response level: with $\mathbf{Z} \in \mathbb{R}^{T \times D}$ (rows are step-wise hidden states along one trace), k measures how many reasoning steps add non-redundant semantic directions, so linear ER/ERV indicates sustained novelty and $\mathcal{O}(1)$ ERA indicates non-saturating refinement.

4 Empirical Evidence for Decoupling in Semantic Space

We empirically study how exploration and exploitation evolve *in the hidden-state semantic space* during RLVR training. Concretely, we track both conventional rank and our semantic-space metrics (ER, ERV, ERA) throughout GRPO post-training. We experiment with Qwen (Hui et al., 2024) and Llama (Dubey et al., 2024) under the GRPO paradigm (Shao et al., 2024), using the training configuration of Zeng et al. (2025): 8k hard MATH problems (level 3-5), each paired with a verifiable reference answer. Complementary theoretical analysis is available in App. E.1.

4.1 Analysis of Response-Level Metrics

At each training step, we compute response-level metrics on the current batch (Overall / Correct / Incorrect) as shown in Fig. 2, with additional analyses deferred to App. G.1.

ER varies by model, while ERV increases consistently across models. Although RL consistently improves task performance, it induces different exploration patterns across base models in semantic space. As shown in the first column of Fig. 2, Qwen exhibits an increasing ER, suggesting broader coverage over semantic directions, whereas Llama shows a decreasing ER, consistent with a more concentrated trajectory. In contrast, ERV (second column) exhibits a consistent upward trend for both models, indicating that RL fine-tuning re-

liably strengthens *semantic refinement dynamics*. Concretely, representations evolve in a way that increasingly outpaces their historical trend, even when the models differ in their exploration breadth.

ERA distinguishes productive refinement from unproductive complexity. We further compare Correct vs. Incorrect trajectories in Fig. 2. Interestingly, Incorrect samples often exhibit higher ER and ERV than Correct ones, suggesting that larger semantic breadth (high ER) or faster growth relative to the past (high ERV) can reflect *unproductive expansion* that drifts away from a correct solution. In contrast, ERA is consistently higher for Correct trajectories, indicating that correct reasoning is characterized not merely by being “large” or “fast”, but by an *increasingly coherent refinement* where the refinement dynamics strengthen over time. This makes ERA a more reliable indicator of robust reasoning than zero- or first-order quantities.

4.2 Analysis of Dataset-level Metrics

Following Sec. 2, we extend the analysis from individual trajectories to the validation set by constructing the dataset hidden-state matrix $\bar{\mathbf{Z}}$. We then track the zero-, first-, and second-order dynamics of $\text{erank}(\bar{\mathbf{Z}})$ during training (Fig. 3), with additional studies in App. G.2.

Policy improvement correlates with expanding dataset-level semantic diversity. Across training, we observe that validation performance (accuracy and score) co-evolves with dataset-level ER dynamics. As optimization proceeds, $\text{erank}(\bar{\mathbf{Z}})$ and its temporal differences tend to increase, suggesting that the learned policy occupies a richer set of semantic directions across the same problem set. The upward trends of ERV and ERA further indi-

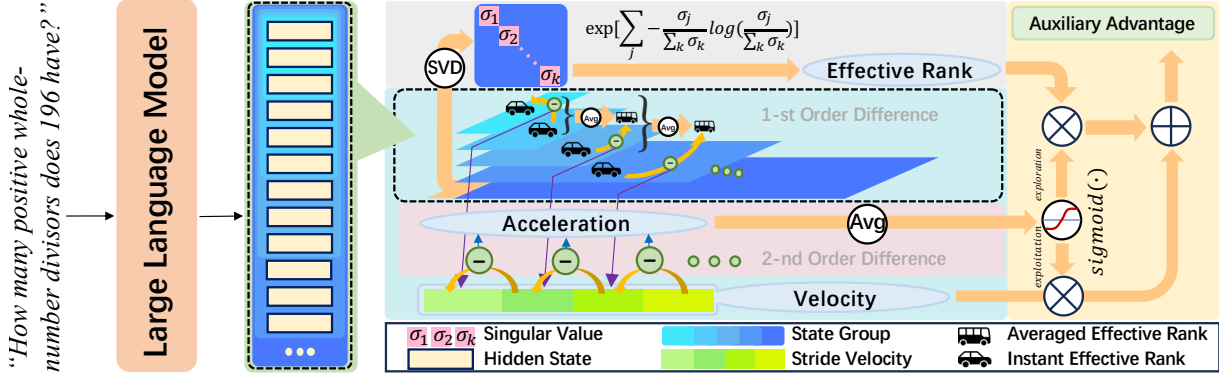


Figure 4: Overview of VERL. Exploration is quantified by computing the ER of the rolling-done hidden states via SVD, while exploitation is captured through EMA-smoothed first-order difference (ERV) on per-step rolling hidden state and extended to second-order difference (ERA). Finally, exploration and exploitation are adaptively integrated to derive the auxiliary advantage.

cate that the policy becomes progressively more effective at refining and expanding this dataset-level semantic space toward correct solutions.

ER reveals refinement even when conventional rank plateaus. In late-stage training, conventional rank may plateau, suggesting that the number of linearly independent directions in $\bar{\mathbf{Z}}$ no longer increases. However, ER can continue to rise, revealing a subtler improvement: the singular spectrum becomes more balanced, meaning the model utilizes existing semantic directions more uniformly rather than concentrating energy on a few dominant ones. This reduced redundancy indicates refinement within the established semantic space, beyond what discrete rank alone can capture.

5 Velocity-Exploiting Rank-Learning

Building on Sec. 2.3 and the empirical findings in Sec. 4, we propose VERL (Velocity-Exploiting Rank Learning), an advantage-shaping approach that injects *semantic-space* signals derived from hidden-state dynamics into standard RL objectives. Conventional RLVR optimizes rewards in the token-level action space, but does not explicitly regulate how reasoning trajectories evolve in representation space. This can yield inefficient updates, e.g., unproductive semantic expansion or premature collapse into narrow reasoning patterns. VERL addresses this by constructing an auxiliary advantage from ER, ERV, and ERA, encouraging both broader semantic coverage and productive refinement in a unified and stable manner.

5.1 Normalized Deviations in Semantic Space

We compute three scalar metrics per trajectory, $\{m_0, m_1, m_2\}$, corresponding to ER (breadth), ERV (first-order refinement dynamics), and ERA

(second-order refinement dynamics). To obtain a stable, scale-invariant control signal across training, we normalize each metric by its recent trend using an exponential moving average (EMA) $\bar{\mu}_k$. We define the relative deviation

$$d_k := (m_k - \bar{\mu}_k) / (|\bar{\mu}_k| + \epsilon), \quad k \in \{0, 1, 2\}, \quad (10)$$

where ϵ is a small constant. Intuitively, d_k measures how the current trajectory deviates from the policy’s recent behavior in semantic space.

5.2 Adaptive Balancing via ERA as a Meta-Control Signal

Sec. 3.1–3.2 motivates ER (M_0) as a proxy for exploration breadth in semantic space and ERV (M_1) as a proxy for refinement dynamics (exploitation). Empirically, these two quantities are nearly uncorrelated (Fig. 1b), suggesting they can be jointly optimized rather than traded off. We therefore combine M_0 and M_1 into a single auxiliary signal, and use ERA (M_2) as a meta-control variable to adapt their relative emphasis.

The key intuition is that ERA captures the *change* of refinement dynamics. When refinement accelerates strongly relative to its EMA trend ($d_2 \gg 0$), the policy may be at risk of collapsing into overly confident, narrow patterns that generalize poorly. When refinement decelerates or saturates ($d_2 \leq 0$), strengthening exploitation becomes more important to consolidate useful progress. This motivates an ERA-controlled interpolation between exploration- and exploitation-oriented weights.

We define two basis weight vectors, $\mathbf{w}_{\text{explore}} = [1, 0]$ for M_0 and $\mathbf{w}_{\text{exploit}} = [0, 1]$ for M_1 . The dynamic weight is

$$\mathbf{w}_{\text{dyn}} := \beta \mathbf{w}_{\text{explore}} + (1 - \beta) \mathbf{w}_{\text{exploit}}, \quad (11)$$

Table 1: Performance comparison of models on mathematical reasoning benchmarks (Pass@1). “+ GRPO” and “+ PPO” denote RL fine-tuning by GRPO and PPO framework respectively. “w/ VERL.” indicates incorporating our VERL with original RL type. Δ represents the performance contrast between original RL method and its VERL variant. See App. G.3 for full details.

Model	AIME24	AIME25	AMC23	AMC24	ASDiv	Carp_En	CMATH	Gaokao 2024_I	Gaokao 2024_Mix	Gaokao MathCloze	GSM8K	MAWPS	Olympiad Bench	SVAMP	TabMWP	Avg.
Llama-3.2-3B-Instruct	0.0	0.0	25.0	11.1	74.6	26.5	10.2	14.3	14.3	6.8	66.6	86.9	12.7	74.1	41.4	31.0
+ GRPO	3.3	0.0	27.5	8.9	88.8	45.0	28.3	21.4	20.9	23.7	80.7	96.0	16.7	87.7	71.7	41.4
+ GRPO w/ VERL.	13.3	6.7	25.0	11.1	89.3	45.4	46.2	14.3	22.0	22.9	81.7	96.0	17.6	87.8	72.3	43.4
Δ_{GRPO}	+10.0	+6.7	-2.5	+2.2	+0.5	+0.4	+17.9	-7.1	+1.1	-0.8	+1.0	+0.0	+0.9	+0.1	+0.6	+2.0
+ PPO	10.0	3.3	22.5	13.3	87.9	46.4	21.2	7.1	16.5	20.3	81.4	95.5	17.8	86.8	71.0	40.1
+ PPO w/ VERL.	10.0	3.3	25.0	11.1	88.7	46.0	30.7	14.3	19.8	27.1	82.9	95.7	17.3	85.8	71.3	41.9
Δ_{PPO}	+0.0	+0.0	+2.5	-2.2	+0.8	-0.4	+9.5	+7.2	+3.3	+6.8	+1.5	+0.2	-0.5	-1.0	+0.3	+1.9
Qwen2.5-7B	6.7	0.0	45.0	15.6	91.4	55.8	86.7	42.9	33.0	49.2	85.8	95.4	25.8	88.5	82.8	53.6
+ GRPO	10.0	6.7	55.0	26.7	94.8	60.2	91.7	14.3	34.1	64.4	90.2	97.6	36.1	92.8	91.3	57.7
+ GRPO w/ VERL.	13.3	10.0	50.0	28.9	95.0	60.8	90.7	35.7	35.2	69.5	89.2	97.7	35.4	92.9	91.9	59.8
Δ_{GRPO}	+3.3	+3.3	-5.0	+2.2	+0.2	+0.6	-1.0	+21.4	+1.1	+5.1	-1.0	+0.1	-0.7	+0.1	+0.6	+2.1
+ PPO	6.7	3.3	50.0	33.3	94.9	59.6	89.8	28.6	31.9	63.6	89.1	97.3	36.1	92.8	90.8	57.9
+ PPO w/ VERL.	10.0	6.7	52.5	33.3	94.8	60.0	90.3	28.6	34.1	66.9	90.2	97.8	36.1	92.5	90.6	59.0
Δ_{PPO}	+3.3	+3.3	+2.5	+0.0	-0.1	+0.4	+0.5	+0.0	+2.2	+3.3	+1.1	+0.5	+0.0	-0.3	-0.2	+1.1

Table 2: Performance comparison of instruction-tuned models under diverse decoding settings (Pass@k). For full details, please refer to App. G.4.

Model	MATH500 (Pass@16)	AMC23 (Pass@128)	AMC24 (Pass@128)	AIME24 (Pass@256)	AIME25 (Pass@256)	Avg.
Llama-3.2-3B-Instruct	79.8	93.5	51.5	40.0	30.0	58.96
+ GRPO	80.2	95.4	60.6	40.0	30.0	61.24
+ GRPO w/ VERL.	80.6	95.7	59.0	50.0	36.7	64.40
Δ_{GRPO}	+0.4	+0.3	-1.6	+10.0	+6.7	+3.16
+ PPO	82.2	94.5	57.0	46.7	36.7	63.42
+ PPO w/ VERL.	82.4	94.7	57.8	46.7	40.0	64.32
Δ_{PPO}	+0.2	+0.2	+0.8	+0.0	+3.3	+0.90
Qwen2.5-7B	90.6	98.4	73.7	60.0	60.0	76.54
+ GRPO	90.8	97.8	78.3	56.7	50.0	74.72
+ GRPO w/ VERL.	91.4	98.3	79.0	63.3	60.0	78.40
Δ_{GRPO}	+0.6	+0.5	+0.7	+6.6	+10.0	+3.68
+ PPO	91.2	98.6	74.3	53.3	56.7	74.82
+ PPO w/ VERL.	91.4	98.0	74.4	56.7	66.7	77.44
Δ_{PPO}	+0.2	-0.6	+0.1	+3.4	+10.0	+2.62

where $\beta := \sigma(d_2)$ and $\sigma(\cdot)$ is the sigmoid. We then define the auxiliary trajectory-level signal

$$\Phi := w_{\text{dyn},0} \cdot \tanh(d_0) + w_{\text{dyn},1} \cdot \tanh(d_1), \quad (12)$$

where \tanh bounds the magnitude for stability. Intuitively, Φ rewards trajectories that exceed recent semantic-space trends in breadth or refinement, while adaptively preventing over-collapse and stagnation via d_2 .

5.3 Advantage Shaping with Auxiliary Semantic-Space Signals

We shape the learning signal by adding a clipped, non-negative bonus derived from Φ . Let A_t denote the original advantage (e.g., from GRPO or PPO with GAE), and let Φ_i be the trajectory-level auxiliary signal in Sec. 5.2. We use the refined advantage

$$\hat{A}_t = A_t + \min(\text{ReLU}(\Phi_i), |A_t|/\kappa), \quad (13)$$

which preserves the sign of A_t while bounding the shaping magnitude for stability. For GRPO, Φ_i is applied once per trajectory; for PPO+GAE, the same Φ_i is broadcast to all tokens in the trajectory.

6 Experiments

6.1 Experiment Settings

(i) **Datasets.** We use the same training and evaluation datasets as in Sec. 4. (ii) **Reward.** We

employ a rule-based reward that checks mathematical correctness and `\boxed{\}` formatting. A correct answer receives +1.0 if properly formatted and +0.5 otherwise; an incorrect answer receives -0.5 if formatted and -1.0 otherwise. (iii) **Training.** We build on VERL (Sheng et al., 2025) with vLLM (Kwon et al., 2023). We use batch size 48, generate 4 rollouts per prompt for GRPO and 1 rollout for PPO, and set the maximum generation length $L_{\text{max}} = 1536$.

6.2 Main Results

Generalization across benchmarks. Tab. 1 summarizes Pass@1 results on 15 mathematical reasoning benchmarks (full tables in App. G.3). Across both easy (e.g., ASDiv) and hard benchmarks (e.g., OlympiadBench), VERL consistently improves over the corresponding RL baselines. Gains are most pronounced on multi-step reasoning benchmarks, with up to **+21.4** and **+10.0** absolute accuracy on Gaokao 2024_I and AIME24 with Qwen2.5-7B, respectively.

Generalization across RL algorithms and base models. VERL is plug-and-play: it can be paired with both GRPO and PPO, and improves average performance for both Llama and Qwen models in Tab. 1. This indicates that semantic-space advantage shaping transfers across optimization

Table 3: Pass@1 performance with variant β . “Adapted β ” denotes $\beta := \text{sigmoid}(d_2)$. In this paper, all results in the table are reported in percentage (%), with **Bold** indicating the best performance.

Training Strategy	Score Avg	In Domain			Out of Domain				Hard Problems				
		MATH	MATH500	Avg	Gaokao	CN Middle School	CMATH	Avg	AIME24	AIME25	AMC23	AMC24	Avg
GRPO	0.36	51.4	46.2	48.80	23.7	28.7	28.3	26.90	3.3	0.0	27.5	8.9	9.93
GRPO+VERL ($\beta = 0.5$)	0.38	51.2	47.2	49.20	21.2	36.6	38.7	32.17	10.0	0.0	27.5	8.9	11.60
GRPO+VERL (Adapted β)	0.38	50.9	51.2	51.05	22.9	38.6	46.2	35.90	13.3	6.7	20.0	11.1	12.78

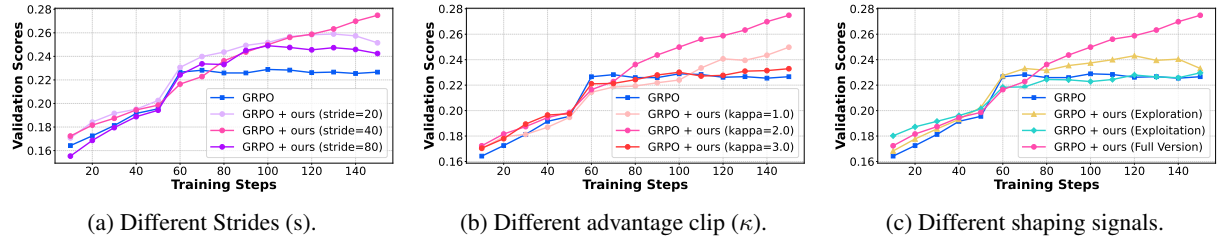


Figure 5: Comparison of various hyperparameters with Llama-3.2-3B-Instruct. It shows that the model performs best with a stride of 40 in (a) and with $\kappa = 2$ in (b). We adopt these settings for all subsequent experiments. Moreover, (c) indicates that using only one signal, either exploration or exploitation, leads to suboptimal performance, demonstrating the effectiveness of our method.

algorithms and base-model families.

Improving both exploration and exploitation proxies. Tab. 2 reports Pass@ k under diverse decoding budgets (App. G.4). Pass@1 is commonly used as a proxy for strong single-shot performance, whereas Pass@ k reflects the ability to succeed given a larger search budget. We observe larger gains on Pass@ k than on Pass@1, especially on harder benchmarks, suggesting that VERL improves not only single-trajectory refinement but also the coverage of successful reasoning trajectories under increased sampling. See Sec. H for qualitative analysis.

6.3 Ablations and Analysis

We ablate key components and hyperparameters of VERL, including the role of ERA, the stride s for temporal metrics, the advantage clip κ , and the composition of the shaping signal Φ .

Effectiveness of ERA. Tab. 3 shows that adding semantic-space shaping improves over GRPO. Moreover, adaptive mixing ($\beta = \sigma(d_2)$) outperforms fixed mixing ($\beta = 0.5$), indicating that ERA provides a useful meta-control signal for adjusting when to emphasize breadth (ER) versus refinement (ERV). Thus, VERL improved the performance on OOD and challenging tasks.

Stride s . The stride s controls the temporal granularity for computing ERV and ERA. As shown in Fig. 5(a), improvements are robust across a range of s , suggesting that the signal is not overly sensitive to sampling frequency. We use $s = 40$ in subsequent experiments as it yields the best val-

idation reward while avoiding noisy token-level fluctuations.

Advantage clip κ . The clip factor κ bounds the shaping bonus so it complements, rather than dominates, the original advantage. As shown in Fig. 5(b), VERL consistently improves performance across all tested κ , indicating robust behavior to this hyperparameter. The best performance is achieved at $\kappa = 2$, which provides a sufficiently strong yet well-proportioned shaping signal.

Shaping signal Φ . Fig. 5(c) compares variants of Eq. 12. Using only the exploration term improves coverage but under-utilizes high-reward trajectories, leading to earlier bottlenecks. Using only the exploitation term yields faster initial gains but plateaus due to insufficient breadth. Combining both terms achieves the most stable training and the best final performance, consistent with our goal of improving both capacities in semantic space.

7 Conclusion

We argue that the exploration–exploitation trade-off in RLVR for LLM reasoning is largely a measurement artifact of token-level action statistics. Moving to hidden-state semantic space, we propose ER, ERV, and ERA to capture semantic breadth and refinement dynamics, and empirically find near-zero correlation between exploration and exploitation. Leveraging ERA as a meta-control signal, VERL performs semantic-space advantage shaping and consistently improves both in-domain performance and out-of-domain generalization on challenging reasoning benchmarks.

561 Limitations

562 We have not systematically studied how training-
563 data cleaning or curation affects VERL. While
564 VERL shows clear advantages on math-focused
565 training data with verifiable answers, extending
566 these gains to other reasoning training corpora re-
567 mains an open direction for further development
568 and investigation.

569 References

570 Shivam Agarwal, Zimin Zhang, Lifan Yuan, Jiawei Han,
571 and Hao Peng. 2025. [The unreasonable effective-
572 ness of entropy minimization in llm reasoning](#). In
573 *Advances in Neural Information Processing Systems*
574 (*NeurIPS*). NeurIPS 2025 Poster.

575 Shelly Bensal, Umar Jamil, Christopher Bryant,
576 Melisa Russak, Kiran Kamble, Dmytro Mozolevskyi,
577 Muayad Ali, and Waseem AlShikh. 2025. Reflect,
578 retry, reward: Self-improving llms via reinforcement
579 learning. *arXiv preprint arXiv:2505.24726*.

580 Zhipeng Chen, Xiaobo Qin, Youbin Wu, Yue Ling,
581 Qinghao Ye, Wayne Xin Zhao, and Guang Shi. 2025.
582 Pass@k training for adaptively balancing exploration
583 and exploitation of large reasoning models. *arXiv*
584 *preprint arXiv:2508.10751*.

585 Daixuan Cheng, Shaohan Huang, Xuekai Zhu, Bo Dai,
586 Wayne Xin Zhao, Zhenliang Zhang, and Furu Wei.
587 2025a. Reasoning with exploration: An entropy per-
588 spective. *arXiv preprint arXiv:2506.14758*.

589 Zhengxiang Cheng, Dongping Chen, Mingyang Fu, and
590 Tianyi Zhou. 2025b. [Lc-rl: Optimizing length com-
591 pression in large reasoning model](#). In *ICML 2025*
592 *Affinity Workshop: New In ML*. Poster.

593 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
594 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
595 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
596 Nakano, and 1 others. 2021. Training verifiers
597 to solve math word problems. *arXiv preprint*
598 *arXiv:2110.14168*.

599 Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan
600 Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen
601 Fan, Huayu Chen, Weize Chen, and 1 others. 2025.
602 The entropy mechanism of reinforcement learning
603 for reasoning language models. *arXiv preprint*
604 *arXiv:2505.22617*.

605 Mehul Damani, Isha Puri, Stewart Slocum, Idan Shen-
606 feld, Leshem Choshen, Yoon Kim, and Jacob An-
607 dreas. 2025. Beyond binary rewards: Training lms
608 to reason about their uncertainty. *arXiv preprint*
609 *arXiv:2507.16806*.

610 Jia Deng, Jie Chen, Zhipeng Chen, Daixuan Cheng,
611 Fei Bai, Beichen Zhang, Yinqian Min, Yanzipeng
612 Gao, Wayne Xin Zhao, and Ji-Rong Wen. 2025a.
613 From trial-and-error to improvement: A systematic

analysis of llm exploration mechanisms in rlvr. *arXiv*
preprint arXiv:2508.07534. 614 615

Wenlong Deng, Yi Ren, Danica J. Sutherland, Christos
Thrampoulidis, and Xiaoxiao Li. 2025b. [Token hid-
den reward: Steering exploration-exploitation in gpo
training](#). In *ICML 2025 Workshop on AI for Math*
(*AI4MATH*). Oral. 616 617 618 619 620

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
Akhil Mathur, Alan Schelten, Amy Yang, Angela
Fan, and 1 others. 2024. The llama 3 herd of models.
arXiv e-prints, pages arXiv–2407. 621 622 623 624 625

Yu Fan, Jingwei Ni, Jakob Merane, Yang Tian, Yoan
Hermstrüwer, Yinya Huang, Mubashara Akhtar, Eti-
enne Salimbeni, Florian Geering, Oliver Dreyer,
and 1 others. 2025. Lexam: Benchmarking legal
reasoning on 340 law exams. *arXiv preprint*
arXiv:2505.12864. 626 627 628 629 630 631

Yichao Fu, Xuewei Wang, Yuandong Tian, and Jiawei
Zhao. 2025. [Deep think with confidence](#). In *NeurIPS*
2025 Workshop on ER. Spotlight. 632 633 634

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song,
Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma,
Peiyi Wang, Xiao Bi, and 1 others. 2025. [Deepseek-
r1 incentivizes reasoning in llms through reinforce-
ment learning](#). *Nature*, 645:633–638. 635 636 637 638 639

Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding
Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han,
Yujie Huang, Yuxiang Zhang, and 1 others. 2024.
Olympiadbench: A challenging benchmark for pro-
moting agi with olympiad-level bilingual multimodal
scientific problems. In *Proceedings of the 62nd An-
nual Meeting of the Association for Computational*
Linguistics (Volume 1: Long Papers), pages 3828–
3850. 640 641 642 643 644 645 646 647 648

Dan Hendrycks, Collin Burns, Steven Basart, Andy
Zou, Mantas Mazeika, Dawn Song, and Jacob Stein-
hardt. 2021a. [Measuring massive multitask language](#)
[understanding](#). In *9th International Conference on*
Learning Representations, ICLR 2021, Virtual Event,
Austria, May 3-7, 2021. OpenReview.net. 649 650 651 652 653 654

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul
Arora, Steven Basart, Eric Tang, Dawn Song, and
Jacob Steinhardt. 2021b. [Measuring mathematical](#)
[problem solving with the math dataset](#). In *Proceed-*
ings of the Neural Information Processing Systems
Track on Datasets and Benchmarks. 655 656 657 658 659 660

Fanding Huang, Jingyan Jiang, Qinting Jiang, Hebei
Li, Faisal Nadeem Khan, and Zhi Wang. 2025.
[COSMIC: clique-oriented semantic multi-space in-
tegration for robust CLIP test-time adaptation](#). In
*IEEE/CVF Conference on Computer Vision and Pat-
tern Recognition, CVPR 2025, Nashville, TN, USA,*
June 11-15, 2025, pages 9772–9781. Computer Vi-
sion Foundation / IEEE. 661 662 663 664 665 666 667 668

669	Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang,	Timothée Lesort, Natalia Díaz Rodríguez, Jean-François	724
670	Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun	Goudou, and David Filliat. 2018. State representation learning for control: An overview . <i>Neural networks : the official journal of the International Neural Network Society</i> , 108:379–392.	725
671	Zhang, Bowen Yu, Keming Lu, and 1 others. 2024.		726
672	Qwen2. 5-coder technical report. <i>arXiv preprint arXiv:2409.12186</i> .		727
673			728
674	Albert Q Jiang, Alexandre Sablayrolles, Antoine	Haitao Li, You Chen, Qingyao Ai, Yueyue Wu, Ruizhe	729
675	Roux, Arthur Mensch, Blanche Savary, Chris Bam-	Zhang, and Yiqun Liu. 2024. Lexeval: A compre-	730
676	ford, Devendra Singh Chaplot, Diego de las Casas,	hensive chinese legal benchmark for evaluating large	731
677	Emma Bou Hanna, Florian Bressand, and 1 oth-	language models. <i>Advances in Neural Information</i>	732
678	ers. 2024. Mixtral of experts. <i>arXiv preprint</i>	<i>Processing Systems</i> , 37:25061–25094.	733
679	<i>arXiv:2401.04088</i> .		
680	Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon,	Jie Li, Fuyong Zhao, Panfeng Chen, Jiafu Xie, Xiangrui	734
681	Sercan O. Arik, Dong Wang, Hamed Zamani, and	Zhang, Hui Li, Mei Chen, Yanhao Wang, and Ming	735
682	Jiawei Han. 2025. Search-r1: Training llms to reason and leverage search engines with reinforcement learning . In <i>Conference on Language Modeling (COLM)</i> .	Zhu. 2025a. An astronomical question answering	736
683		dataset for evaluating large language models. <i>Scientific Data</i> , 12(1):447.	737
684			738
685	Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng,	Pengyi Li, Matvey Skripkin, Alexander Zubrey, Andrey	739
686	Hanyi Fang, and Peter Szolovits. 2021. What disease does this patient have? a large-scale open domain question answering dataset from medical exams . <i>Applied Sciences</i> , 11(14):6421.	Kuznetsov, and Ivan Oseledets. 2025b. Confidence	740
687		is all you need: Few-shot rl fine-tuning of language	741
688		models. <i>arXiv preprint arXiv:2506.06395</i> .	742
689			
690	Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W.	Xiao Liang, Zhong-Zhi Li, Yeyun Gong, Yang Wang,	743
691	Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 2567–2577. Association for Computational Linguistics.	Hengyuan Zhang, Yelong Shen, Ying Nian Wu, and	744
692		Weizhu Chen. 2025a. Sws: Self-aware weakness-driven problem synthesis in reinforcement learning for llm reasoning . In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> . NeurIPS 2025 Poster.	745
693			746
694			747
695			748
696		Xiao Liang, Zhongzhi Li, Yeyun Gong, Yelong Shen,	749
697		Ying Nian Wu, Zhijiang Guo, and Weizhu Chen.	750
698		2025b. Beyond pass@ 1: Self-play with variation-	751
699		al problem synthesis sustains rlvr. <i>arXiv preprint</i>	752
700		<i>arXiv:2508.14029</i> .	753
701	Yi Jing, Zijun Yao, Hongzhu Guo, Lingxu Ran, Xiaozhi	Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang,	754
702	Wang, Lei Hou, and Juanzi Li. 2025. Lingualens: Towards interpreting linguistic mechanisms of large language models via sparse auto-encoder . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 28220–28239.	Yile Wang, and Yue Zhang. 2020. Logiqa: A challenge dataset for machine reading comprehension with logical reasoning . In <i>Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20</i> , pages 3622–3628. International Joint Conferences on Artificial Intelligence Organization. Main track.	755
703			756
704			757
705	Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate	Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu,	758
706	Kushman, and Hannaneh Hajishirzi. 2016. Mawps: A math word problem repository . In <i>Proceedings of the 2016 conference of the north american chapter of the association for computational linguistics: human language technologies</i> , pages 1152–1157.	Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark,	759
707		and Ashwin Kalyan. 2023. Dynamic prompt learning	760
708		via policy gradient for semi-structured mathematical	761
709		reasoning. In <i>ICLR</i> .	762
710			763
711	Aviral Kumar, Rishabh Agarwal, Dibya Ghosh,	Clare Lyle, Mark Rowland, and Will Dabney. 2022.	764
712	and Sergey Levine. 2021. Implicit underparameterization inhibits data-efficient deep reinforcement learning . In <i>Proceedings of the 9th International Conference on Learning Representations (ICLR 2021)</i> .	Understanding and preventing capacity loss in reinforcement learning . In <i>Proceedings of the 10th International Conference on Learning Representations (ICLR 2022)</i> .	765
713			766
714			767
715			768
716			769
717	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	Jacob A. Matthews, John R. Starr, and Marten van Schi-	770
718	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gon-	ndel. 2024. Semantics or spelling? probing contextual word embeddings with orthographic noise . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 4445–4462, Bangkok, Thailand. Association for Computational Linguistics.	771
719	zalez, Hao Zhang, and Ion Stoica. 2023. Efficient		772
720	memory management for large language model serv-		773
721	ing with pagedattention. In <i>Proceedings of the 29th symposium on operating systems principles</i> , pages 611–626.		774
722			775
723			776
		Shen-Yun Miao, Chao-Chun Liang, and Keh-Yih Su.	777
		2020. A diverse corpus for evaluating and developing	778
			779

780	english math word problem solvers. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 975–984.	
781		
782		
783	Ankit Pal, Logesh Kumar Umapathi, and Malaikanan Sankarasubbu. 2022. Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering . In <i>Conference on Health, Inference, and Learning, CHIL 2022, 7-8 April 2022, Virtual Event</i> , volume 174 of <i>Proceedings of Machine Learning Research</i> , pages 248–260. PMLR.	
784		
785		
786		
787		
788		
789		
790	Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are nlp models really able to solve simple math word problems? In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2080–2094.	
791		
792		
793		
794		
795		
796	Ziqing Qiao, Yongheng Deng, Jiali Zeng, Dong Wang, Lai Wei, Guanbo Wang, Fandong Meng, Jie Zhou, Ju Ren, and Yaoxue Zhang. 2025. ConCISE: Confidence-guided compression in step-by-step efficient reasoning . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 8021–8040, Suzhou, China. Association for Computational Linguistics.	
797		
798		
799		
800		
801		
802		
803		
804	Corrado Rainone, Tim Bakker, and Roland Memisevic. 2025. Replacing thinking with tool usage enables reasoning in small language models . In <i>ICML 2025 Workshop on Test-Time Adaptation: Putting Updates to the Test (PUT)</i> . Poster.	
805		
806		
807		
808		
809	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2024. GPQA: A graduate-level google-proof q&a benchmark . In <i>First Conference on Language Modeling</i> .	
810		
811		
812		
813		
814	Olivier Roy and Martin Vetterli. 2007. The effective rank: A measure of effective dimensionality. In <i>2007 15th European signal processing conference</i> , pages 606–610. IEEE.	
815		
816		
817		
818	Hassan Sajjad, Nadir Durrani, Fahim Dalvi, Firoj Alam, Abdul Khan, and Jia Xu. 2022. Analyzing encoded concepts in transformer language models. In <i>Proceedings of the 2022 Conference of the North American chapter of the Association for Computational Linguistics: Human language technologies</i> , pages 3082–3101.	
819		
820		
821		
822		
823		
824		
825	John Schulman, Philipp Moritz, Sergey Levine, Michael I. Jordan, and Pieter Abbeel. 2016. High-dimensional continuous control using generalized advantage estimation . In <i>4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings</i> .	
826		
827		
828		
829		
830		
831		
832	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> .	
833		
834		
835		
	Giovanni Servedio, Alessandro De Bellis, Dario Di Palma, Vito Walter Anelli, and Tommaso Di Noia. 2025. Are the hidden states hiding something? testing the limits of factuality-encoding capabilities in llms . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025</i> , pages 6089–6104. Association for Computational Linguistics.	836 837 838 839 840 841 842 843 844
	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y.K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models .	845 846 847 848 849
	Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2025. Hybridflow: A flexible and efficient rlhf framework. In <i>Proceedings of the Twentieth European Conference on Computer Systems</i> , pages 1279–1297.	850 851 852 853 854 855
	Oscar Skean, Md Rifat Arefin, Dan Zhao, Niket Patel, Jalal Naghiyev, Yann LeCun, and Ravid Shwartz-Ziv. 2025. Layer by layer: Uncovering hidden representations in language models . In <i>Forty-second International Conference on Machine Learning, ICML 2025, Vancouver, BC, Canada, July 13-19, 2025</i> . OpenReview.net.	856 857 858 859 860 861 862
	Lucrezia Valeriani, Diego Doimo, Francesca Cuturello, Alessandro Laio, Alessio Ansuini, and Alberto Cazzaniga. 2023. The geometry of hidden representations of large transformer models. <i>Advances in Neural Information Processing Systems</i> , 36:51234–51252.	863 864 865 866 867 868
	Han Wang, Erfan Miah, Martha White, Marlos C. Machado, Zaheer Abbas, Raksha Kumaraswamy, Vincent Liu, and Adam White. 2024a. Investigating the properties of neural network representations in reinforcement learning . <i>Artificial Intelligence</i> , page 104100.	869 870 871 872 873 874
	Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xiong-Hui Chen, Jianxin Yang, Zhenru Zhang, Yuqiong Liu, An Yang, Andrew Zhao, Yang Yue, Shiji Song, Bowen Yu, Gao Huang, and Junyang Lin. 2025. Beyond the 80/20 rule: High-entropy minority tokens drive effective reinforcement learning for llm reasoning . In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> . NeurIPS 2025 Poster.	875 876 877 878 879 880 881 882 883
	Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R. Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. 2024b. SciBench: Evaluating College-Level Scientific Problem-Solving Abilities of Large Language Models. In <i>Proceedings of the Forty-First International Conference on Machine Learning</i> .	884 885 886 887 888 889 890
	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren,	891 892

893	Aaran Arulraj, Xuan He, Ziyang Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhui Chen. 2024c. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	Anqi Zhang, Yulin Chen, Jane Pan, Chen Zhao, Aurojit Panda, Jinyang Li, and He He. 2025. Reasoning models know when they’re right: Probing hidden states for self-verification . In <i>Second Conference on Language Modeling (COLM)</i> .	951 952 953 954 955
901	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	Beichen Zhang, Kun Zhou, Xilin Wei, Xin Zhao, Jing Sha, Shijin Wang, and Ji-Rong Wen. 2023. Evaluating and improving tool-augmented computation-intensive math reasoning. <i>Advances in Neural Information Processing Systems</i> , 36:23570–23589.	956 957 958 959 960
907	Tianwen Wei, Jian Luan, Wei Liu, Shuang Dong, and Bin Wang. 2023. Cmath: Can your language model pass chinese elementary school math test? <i>arXiv preprint arXiv:2306.16636</i> .	Yuxin Zuo, Kaiyan Zhang, Li Sheng, Shang Qu, Ganqu Cui, Xuekai Zhu, Haozhan Li, Yuchen Zhang, Xinwei Long, Ermo Hua, Biqing Qi, Youbang Sun, Zhiyuan Ma, Lifan Yuan, Ning Ding, and Bowen Zhou. 2025. Ttrl: Test-time reinforcement learning . In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> . NeurIPS 2025 Poster.	961 962 963 964 965 966 967
911	Hang Yan, Fangzhi Xu, Rongman Xu, Yifei Li, Jian Zhang, Haoran Luo, Xiaobao Wu, Luu Anh Tuan, Haiteng Zhao, Qika Lin, and 1 others. 2025. Mur: Momentum uncertainty guided reasoning for large language models . <i>arXiv preprint arXiv:2507.14958</i> .		
916	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, and 1 others. 2024. Qwen2 technical report. <i>arXiv preprint arXiv:2407.10671</i> .		
921	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. <i>Advances in neural information processing systems</i> , 36:11809–11822.		
926	Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gao-hong Liu, Juncai Liu, LingJun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, and 2 others. 2025. Dapo: An open-source llm reinforcement learning system at scale . In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> . NeurIPS 2025 Poster.		
934	Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. 2020. Reclor: A reading comprehension dataset requiring logical reasoning . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.		
940	Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang. 2025. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model? In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> . NeurIPS 2025 Oral.		
946	Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. 2025. Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the wild . <i>Preprint</i> , arXiv:2503.18892.		

968	Contents of Appendix	
969	A Notations	14
970	B Algorithm	15
971	C Related Work	16
972	C.1 Reinforcement Learning with Veri-	
973	fiable Rewards	16
974	C.2 Exploration and Exploitation in	
975	RLVR of LLM	16
976	C.3 Representation Dynamics in Deep	
977	Reinforcement Learning	16
978	D Details of Theorems	17
979	D.1 Proof of Theorem 3.1	17
980	D.2 Proof of Proposition 3.5	17
981	E Additional Theoretical Support for Ex-	
982	ploration and Exploitation Metrics	19
983	E.1 Action vs. Semantic Exploration–	
984	Exploitation Metrics	19
985	E.2 Effective Rank as Semantic Explo-	
986	ration	21
987	E.3 Effective Rank Velocity as Seman-	
988	tic Exploitation	22
989	F Implementation Details	23
990	F.1 Training and Evaluation Details .	23
991	F.2 Efficient Incremental Computation	
992	of Higher-Order Metrics	24
993	F.3 Time Overhead of VERL Training	25
994	G More Experiments	25
995	G.1 Detailed Analysis of Response-	
996	Level Hidden States	25
997	G.2 Detailed Analysis of Dataset-Level	
998	Hidden States	25
999	G.3 Detailed Analysis of Pass@1 Per-	
1000	formance	25
1001	G.4 Detailed Analysis of Pass@ k Per-	
1002	formance	25
1003	G.5 Robustness Across Random Seeds	30
1004	G.6 Detailed Analysis of OOD Perfor-	
1005	mance	30
1006	G.7 Comparison with Regularization	
1007	Method in Action Space	32
1008	G.8 Ablation on the Choice of Hidden	
1009	Layer	33
1010	H Case Study	34
1011	H.1 Case Study For Pass@1 Setting . .	34
1012	H.2 Case Study For Pass@16 Setting .	34

A Notations

Symbol	Description
π	The policy of Large language model
θ	Large language model policy's parameter
ϕ	The parameter corresponding to the optimal policy
t	Time step
T	Output length
L_{\max}	The maximum length of model's output
x	Prompt
\mathcal{P}_x	Distribution of prompts
y_t	The t -th step (token) of the model's response
$y_{i:j}$	Sequence of reasoning steps from i to j
D	Feature dimension of hidden states
N	The size of the dataset
z_t	Hidden state corresponding to the t -th step of the output token
\bar{z}_i	Single vector representation for the i -th response by averaging its token hidden states
\mathbf{Z}_c	Mean-centered hidden state matrix
$\bar{\mathbf{Z}}_{1:n}$	Dataset-level hidden states matrix formed by the first n prompts
\mathbf{Z}	Response-level hidden states matrix
$\text{SVD}(\cdot)$	The function to calculate the singular values
$\lambda_j(\cdot)$	The j -th eigenvalues of the given matrix
σ_j	The j -th singular value of matrix
p_j	The j -th normalized singular values
\mathbf{p}	Singular value distribution
s	The stride for effective rank velocity calculation
$\text{rank}(\cdot)/\text{erank}(\cdot)$	Conventional rank/Effective rank function
$\Delta_M^{(i)}$	The i -order temporal difference for metric M
\mathcal{M}	Set of metrics derived from the hidden states
M_i	The i -order temporal difference of ER, exactly the different metrics
$\bar{\mu}_k$	Exponential moving average of metric M_k over training steps
m_t	Value of metric M computed on the token sub-sequence from the start to position t
$\delta_n^{(i)}$	Instantaneous i -order difference for step n
$\epsilon_{\text{high/low}}$	Hyperparameter for the upper/lower bound used for clipping
ϵ	Small constant for numerical stability
$\mathbf{w}_{\text{explore/exploit/dyn}}$	Exploration-focused profile/Exploitation-focused profile/Dynamic-weighted profile
$w_{\text{dyn},i}$	The i -th scalar of $\mathbf{w}_{\text{dynamic}}$
r_j	Reward of the j -th response
$A_{i,t}$	Group-relative advantage for the t -th token in the i -th response in group
A	The original advantage estimation
\hat{A}	The reshaped advantage value
Φ	Auxiliary advantage
β	Interpolation coefficient for VERL training
G	Size of sampled group in GRPO
ρ_t	Probability ratio between the current and old policies for t -step of the output
$\mathcal{L}_{\text{PPO}}(\cdot)$	The optimization objective for PPO applied to policy
\mathcal{T}	Set of time steps
κ	Advantage clipping factor

B Algorithm

Algorithm 1 VERL: Training

```

1: Input: dataset  $\mathcal{D}$ , prompt  $x \in \mathcal{D}$ ,  $|\mathcal{D}| = N$ , policy model  $\pi_\theta$ , hidden-state dimension  $D$  and sample
   times per prompt  $S$ .
2: Parameters: EMA factor  $\gamma$ , relative deviation stabilizing factor  $\epsilon \ll 1$ , RL fine-tuning stabilizing
   factor  $\kappa$ .
3: Initialize: Randomly initialize policy parameters  $\pi_\theta$ , historical averages of metrics  $\bar{\mu}_{\text{ER}} = \bar{\mu}_{\text{ERV}} =$ 
 $\bar{\mu}_{\text{ERA}} = 0$ , exploration capacity profile  $\mathbf{w}_{\text{explore}} = [1, 0]$ , exploitation capacity profile  $\mathbf{w}_{\text{exploit}} = [0, 1]$ .
4: Output: A well-trained policy model  $\pi_\theta$ .
5: repeat
6:   for  $x \in \mathcal{D}$  do:                                     // Pick a sample from dataset
7:     for 1 to  $S$  do:                                       // Rolling  $S$  times for one sample
8:        $y_0 \leftarrow x, \mathbf{Z}_0 \leftarrow \emptyset, t \leftarrow 1$ 
9:       repeat                                               // Generation process
10:         $y_t, \mathbf{z}_t \sim \pi_\theta(\cdot | y_{t-1})$ 
11:         $y_t \leftarrow [y_{t-1}; y_t]$                                // Concatenate token sequence
12:         $\mathbf{Z}_t \leftarrow [\mathbf{Z}_{t-1}; \mathbf{z}_t] \in \mathbb{R}^{t \times D}$ 
13:         $\sigma_t \leftarrow \text{SVD}(\mathbf{Z}_t)$ 
14:         $j \leftarrow |\sigma_t|$ 
15:         $p_{j,t} \leftarrow \sigma_{j,t} / \sum_j \sigma_{j,t}$ 
16:         $\text{erank}_t \leftarrow \exp\left(-\sum_j p_{j,t} \log p_{j,t}\right)$ 
17:        If  $t > 1$  then:  $\delta_{\text{ERV},t} \leftarrow \text{erank}_t - \frac{1}{t-1} \sum_{k=1}^{t-1} \text{erank}_k$ 
18:        If  $t > 2$  then:  $\delta_{\text{ERA},t} \leftarrow \delta_{\text{ERV},t} - \delta_{\text{ERV},t-1}$ 
19:         $t \leftarrow t + 1$ 
20:      until rolling done the sentence;                       //  $t - 1$  is the final timestep while rolling done
21:       $A_{\text{origin}} \leftarrow$  base RL evaluating on  $y_{t-1}$ 
22:       $m_{\text{ER}} \leftarrow \text{erank}_{t-1}$                                // Calculating ER metric
23:       $m_{\text{ERV}} \leftarrow \frac{1}{t-2} \sum_{t=2}^{t-1} \delta_{\text{ERV}}^t$            // Calculating ERV metric
24:       $m_{\text{ERA}} \leftarrow \frac{1}{t-3} \sum_{t=3}^{t-1} \delta_{\text{ERA}}^t$            // Calculating ERA metric
25:       $\bar{\mu}_k \leftarrow \gamma \bar{\mu}_k + (1 - \gamma) m_k, k \in \{\text{ER}, \text{ERV}, \text{ERA}\}$ 
26:       $d_k \leftarrow \frac{m_k - \bar{\mu}_k}{|\bar{\mu}_k| + \epsilon}, k \in \{\text{ER}, \text{ERV}, \text{ERA}\}$ 
27:       $\beta \leftarrow \text{sigmoid}(d_{\text{ERA}})$ 
28:       $\mathbf{w}_{\text{dyn}} \leftarrow \beta \mathbf{w}_{\text{explore}} + (1 - \beta) \mathbf{w}_{\text{exploit}}$ 
29:       $w_{\text{dyn,ER}} \leftarrow$  the first scalar value of  $\mathbf{w}_{\text{dyn}}$ 
30:       $w_{\text{dyn,ERV}} \leftarrow$  the second scalar value of  $\mathbf{w}_{\text{dyn}}$ 
31:       $\Phi \leftarrow w_{\text{dyn,ER}} \tanh(d_{\text{ER}}) + w_{\text{dyn,ERV}} \tanh(d_{\text{ERV}})$ 
32:       $\hat{A} \leftarrow A_{\text{origin}} + \min\left(\text{ReLU}(\Phi), \frac{|A_{\text{origin}}|}{\kappa}\right)$ 
33:    end for
34:  end for
35:  Update  $\theta$  via base RL objective with  $\hat{A}$ 
36: until  $\theta$  converges;

```

C Related Work

C.1 Reinforcement Learning with Verifiable Rewards

To unlock the full reasoning potential of Large Language Models (LLMs), Reinforcement Learning with Verifiable Rewards (RLVR) has become a prominent training paradigm. This approach was notably employed by DeepSeek-R1-Zero (Guo et al., 2025) and , which executes complex reasoning processes through actions such as reflection and validation. Following the success of DeepSeek-R1, a significant body of research has investigated the efficacy of RLVR on popular open-source LLMs, including Qwen (Yang et al., 2024), Mistral (Jiang et al., 2024), and LLaMA (Dubey et al., 2024).

This has fostered an optimistic view that RLVR can not only enhance existing model capabilities but also enable the acquisition of novel reasoning knowledge, facilitating a path toward continuous self-improvement (Zeng et al., 2025; Yu et al., 2025). RLVR training has been shown to grant LLMs controllable output length for efficient inference (Yan et al., 2025; Cheng et al., 2025b), deepen their reasoning pathways (Bensal et al., 2025), mitigate their weaknesses (Liang et al., 2025a,b), enable the use of external tools (Rainone et al., 2025; Jin et al., 2025), and even facilitate unsupervised reasoning (Zuo et al., 2025). However, some studies (Yue et al., 2025) argue that while RLVR significantly improves the confidence and reliability of model reasoning, it may inadvertently constrain the model’s exploratory capacity. The core of this issue lies in RLVR’s optimization objective: maximizing expected rewards. This objective function inherently biases the policy gradient toward reinforcing known trajectories that lead to high rewards (i.e., “exploitation”), while suppressing the exploration of unknown paths that may offer potentially higher returns but also carry greater risk (i.e., “exploration”). Consequently, the outputs of RLVR-optimized models often remain confined within the sampling distribution of the base model, suggesting the paradigm excels at refining existing knowledge rather than generating new knowledge. This trade-off between exploration and exploitation constitutes a central challenge in the contemporary RLVR landscape.

C.2 Exploration and Exploitation in RLVR of LLM

Recent perspectives (Wang et al., 2025; Cui et al., 2025) on the exploration-exploitation dilemma

have predominantly been shaped by analyses at the token level, focusing on the prediction distribution over the vocabulary. Under this view, next-token entropy serves as an exploration proxy: higher entropy corresponds to a more uniform predictive distribution and is assumed to allow more diverse continuations. This has led to the adoption of techniques (Deng et al., 2025a; Cheng et al., 2025a) such as entropy regularization to explicitly encourage the policy to explore novel reasoning paths. Conversely, lower entropy in the token-level prediction is taken to signify higher model confidence in its reasoning chain, thus representing strong exploitation. Subsequent work (Fu et al., 2025) has also utilized metrics derived from the top-k probabilities of the token prediction to quantify confidence. For instance, some approaches (Damani et al., 2025; Qiao et al., 2025) leverage the model’s internal “confidence” signals to dynamically evaluate and filter the quality of reasoning steps, while others have employed high confidence as a feedback signal to enable unsupervised reinforcement learning (Li et al., 2025b). Ultimately, however, these confidence-based metrics are not fundamentally different from entropy. This token-level standard of measurement introduces an endogenous contradiction: classifying behavior as either exploratory or exploitative requires the introduction of a prior assumptions, a practice that is disadvantageous for LLM research.

In this paper, we depart from this paradigm. We shift the analysis from the token-level action space to the semantic space at the response level. This approach allows us to decouple the intertwined elements of exploration and exploitation, aiming to achieve a simultaneous enhancement of both during reinforcement learning.

C.3 Representation Dynamics in Deep Reinforcement Learning

Beyond RLVR for language models, there is a line of classical deep RL work that explicitly studies how neural representations evolve during training and how this affects exploration and sample efficiency. State representation learning (SRL) for control aims to construct low-dimensional, action-dependent embeddings that preserve task-relevant dynamics while discarding nuisance variation. This work (Lesort et al., 2018) provide a comprehensive overview of SRL methods for robotics and control, emphasizing how compact latent states can improve both data efficiency and stability of

downstream RL algorithms. More recently, several works have directly analyzed the feature dynamics of deep RL agents. The work (Kumar et al., 2021) identify an “implicit under-parameterization” phenomenon in value-based deep RL: repeated bootstrapping updates lead to a collapse in the effective rank of value-network features, which in turn correlates with degraded performance in both on-line and offline settings. The work (Lyle et al., 2022) further study capacity loss, showing that networks trained on non-stationary targets can lose their ability to fit new value functions over time, and proposing Initial Feature Regularization (InFeR) to stabilize the feature subspace and improve performance on sparse-reward Atari tasks. Complementary to these analyses, (Wang et al., 2024a) systematically measure multiple representational properties (e.g., dynamics-awareness, orthogonality) across thousands of agent–task combinations, and relate them to transfer performance, highlighting that not all good control policies arise from equally useful representation geometries.

D Details of Theorems

D.1 Proof of Theorem 3.1

Suppose we have a matrix of embeddings $\mathbf{Z} \in \mathbb{R}^{T \times D}$. Then the effective rank of \mathbf{Z} is a lower bound of $\text{rank}(\mathbf{Z})$:

$$1 \leq \text{erank}(\mathbf{Z}) \leq \text{rank}(\mathbf{Z}) \leq \min\{T, D\} \quad (14)$$

Proof. Let the singular value distribution of the matrix \mathbf{Z} be $\mathbf{p} = (p_1, p_2, \dots, p_{\min\{T, D\}})$. The Shannon entropy of this distribution $H(\mathbf{p})$ is bounded. Its minimum is 0, which occurs when only one element of p is 1 and all others are 0. Its maximum is $\log k$, where k is the number of non-zero singular values, and this occurs when the distribution is uniform ($p_j = 1/k$ for all non-zero values). The lower bound is established from the minimum entropy value:

$$\text{erank}(\mathbf{Z}) = \exp(H(\mathbf{p})) \geq \exp(0) = 1 \quad (15)$$

Equality holds if and only if the singular value distribution is $(1, 0, \dots, 0)$, meaning \mathbf{Z} has only one non-zero singular value. For the upper bound, let $k = \text{rank}(\mathbf{Z})$ be the number of non-zero singular values of \mathbf{Z} . The entropy of the distribution p is calculated only over these k values and is maximized when they are uniform. Therefore

$$H(\mathbf{p}) \leq \log k \quad (16)$$

Applying the exponential function to this inequality gives:

$$\begin{aligned} \text{erank}(\mathbf{Z}) &= \exp(H(\mathbf{p})) \\ &\leq \exp(\log k) = k = \text{rank}(\mathbf{Z}). \end{aligned} \quad (17)$$

This establishes that the effective rank is upper-bounded by the conventional rank. The final inequality, $\text{rank}(\mathbf{Z}) \leq \min\{T, D\}$, is a standard property of matrix rank. Equality for $\text{erank}(\mathbf{Z}) = \text{rank}(\mathbf{Z})$ holds if and only if the non-zero singular values are all equal, corresponding to a uniform singular value distribution over its support.

D.2 Proof of Proposition 3.5

The zero-order metric and first-order difference of the effective rank scales linearly with the number of responses, $\Delta_M^{(0)} = \mathcal{O}(N)$, $\Delta_M^{(1)} = \mathcal{O}(N)$. The second-order difference of the effective rank is constant and does not depend on N , yielding a scaling order of $\Delta_M^{(2)} = \mathcal{O}(1)$.

Proof. Without loss of generality, we take the effective rank for example. We adopt the provided definition of effective rank for a representation matrix \mathbf{Z} with singular values $\{\sigma_i\}$:

$$\text{erank}(\mathbf{Z}) = \exp\left(-\sum_j p_j \log(p_j)\right), \quad (18)$$

$$\text{where } p_j = \frac{\sigma_j}{\sum_k \sigma_k}.$$

Our analysis focuses on the dataset matrix $\bar{\mathbf{Z}} \in \mathbb{R}^{N \times D}$, whose rows $\{\mathbf{q}_i\}_{i=1}^N$ are the mean token embeddings of N responses. The singular values $\sigma_i(\bar{\mathbf{Z}})$ of $\bar{\mathbf{Z}}$ are the square roots of the eigenvalues of the Gram matrix $\bar{\mathbf{K}} = \bar{\mathbf{Z}}\bar{\mathbf{Z}}^\top$; i.e., $\sigma_j(\bar{\mathbf{Z}}) = \sqrt{\lambda_j(\bar{\mathbf{K}})}$. Given that the rows of $\bar{\mathbf{Z}}$ are nearly orthogonal, the Gram matrix $\bar{\mathbf{K}}$ is strongly diagonal-dominant. Its eigenvalues can be approximated by its diagonal entries:

$$\lambda_j(\bar{\mathbf{K}}) \approx \bar{\mathbf{K}}_{jj} = \|\mathbf{q}_j\|^2 = \frac{1}{T} \quad \text{for } j = 1, \dots, N \quad (19)$$

The matrix has N significant eigenvalues, each approximately equal to $1/T$. The singular values of $\bar{\mathbf{Z}}$ are the square roots of the eigenvalues of $\bar{\mathbf{K}}$:

$$\sigma_j(\bar{\mathbf{Z}}) = \sqrt{\lambda_j(\bar{\mathbf{K}})} \approx \sqrt{\frac{1}{T}} = \frac{1}{\sqrt{T}} \quad \text{for } j = 1, \dots, N \quad (20)$$

To calculate the effective rank, we first normalize these singular values to form a probability distribution $\{p_j\}$. The sum of singular values is:

$$\sum_{k=1}^N \sigma_k(\bar{\mathbf{Z}}) \approx \sum_{k=1}^N \frac{1}{\sqrt{T}} = \frac{N}{\sqrt{T}} \quad (21)$$

The individual probabilities are therefore:

$$p_j = \frac{\sigma_j}{\sum_k \sigma_k} \approx \frac{1/\sqrt{T}}{N/\sqrt{T}} = \frac{1}{N} \quad (22)$$

The distribution $\mathbf{p} = \{p_1, p_2, \dots, p_N\}$ is a uniform distribution over N states. The Shannon entropy of this distribution is maximal:

$$\begin{aligned} H(\mathbf{p}) &= -\sum_{j=1}^N p_j \log(p_j) = -\sum_{j=1}^N \frac{1}{N} \log\left(\frac{1}{N}\right) \\ &= -N \left(\frac{-\log(N)}{N}\right) = \log(N). \end{aligned} \quad (23)$$

The effective rank is the exponential of this entropy: $\text{erank}(\bar{\mathbf{Z}}) = \exp(H(\mathbf{p})) = \exp(\log(N)) = N$. In the maximal prompt entropy regime, the effective rank of the dataset matrix scales as $\mathcal{O}(N)$.

We adapt them to our context by defining the metric's value at "time" n as the Effective Rank computed on the dataset matrix formed by the first n prompts, denoted $\bar{\mathbf{Z}}_{1:n}$. Let $m_n = \text{erank}(\bar{\mathbf{Z}}_{1:n})$. From our previous analysis, we established a crucial result that forms the basis of this derivation: for maximal cases, the effective rank of a dataset with n prompts scales linearly with n .

$$m_n = \text{erank}(\bar{\mathbf{Z}}_{1:n}) \approx n \quad (24)$$

We will use this linear approximation to derive the scaling orders of the difference metrics, assuming a stride of $s = 1$ for simplicity. The first-order difference quantifies the average "velocity" of change in the metric relative to its historical mean. Instantaneous First-Order Difference ($\delta_n^{(1)}$) is the value at step n minus the average of all preceding values.

$$\delta_n^{(1)} = m_n - \left(\frac{1}{n-1} \sum_{k=1}^{n-1} m_k\right) \quad (25)$$

Substituting our approximation $m_k \approx k$:

$$\delta_n^{(1)} \approx n - \left(\frac{1}{n-1} \sum_{k=1}^{n-1} k\right) \quad (26)$$

$$\delta_n^{(1)} \approx n - \left(\frac{1}{n-1} \cdot \frac{(n-1)n}{2}\right) = n - \frac{n}{2} = \frac{n}{2} \quad (27)$$

The instantaneous difference grows linearly with n . Overall First-Order Difference $\Delta_{\text{erank}}^{(1)}$: This is the average of the instantaneous differences over the entire dataset of size N .

$$\Delta_{\text{erank}}^{(1)} = \frac{1}{N-1} \sum_{n=2}^N \delta_n^{(1)} \approx \frac{1}{N-1} \sum_{n=2}^N \frac{n}{2} \quad (28)$$

$$\begin{aligned} \Delta_{\text{erank}}^{(1)} &\approx \frac{1}{2(N-1)} \left(\left(\sum_{n=1}^N n \right) - 1 \right) \\ &= \frac{1}{2(N-1)} \left(\frac{N(N+1)}{2} - 1 \right). \end{aligned} \quad (29)$$

For large N , the expression is dominated by the highest power of N :

$$\Delta_{\text{erank}}^{(1)} \sim \frac{N^2/4}{N} = \frac{N}{4} \quad (30)$$

The first-order difference of the effective rank scales linearly with the number of prompts, $\Delta_{\text{erank}}^{(1)} = \mathcal{O}(N)$. As for second-order difference, we compute the change in Instantaneous Differences between consecutive values of $\delta_n^{(1)}$.

$$\delta_n^{(1)} - \delta_{n-1}^{(1)} \approx \frac{n}{2} - \frac{n-1}{2} = \frac{1}{2} \quad (31)$$

This change is a constant, indicating a linear increase in the first-order difference. Overall Second-Order Difference $\Delta_{\text{erank}}^{(2)}$:

$$\begin{aligned} \Delta_{\text{erank}}^{(2)} &= \frac{1}{N-2} \sum_{n=3}^N (d_n^{(1)} - d_{n-1}^{(1)}) \\ &\approx \frac{1}{N-2} \sum_{n=3}^N \frac{1}{2}. \end{aligned} \quad (32)$$

$$\Delta_{\text{erank}}^{(2)} \approx \frac{1}{N-2} \cdot (N-2) \cdot \frac{1}{2} = \frac{1}{2} \quad (33)$$

The second-order difference of the effective rank is constant and does not depend on N , yielding a scaling order of $\Delta_{\text{erank}}^{(2)} = \mathcal{O}(1)$.

E Additional Theoretical Support for Exploration and Exploitation Metrics

In this section we formalize the relationship between our proposed hidden-state metrics (*Effective Rank* and *Effective Rank Velocity*) and the classical notions of exploration and exploitation in reinforcement learning. We first show that the token-level metrics in action space (average log probability and response entropy) are algebraically coupled, whereas the hidden-state metrics in semantic space are not. We then provide a representation-level justification for interpreting *Effective Rank* as a measure of semantic exploration, and *Effective Rank Velocity* as a measure of representation-level exploitation that is strongly correlated with greedy value improvement under the PPO-style architecture used in RLVR. Throughout, we consider a conditional language model $p_\theta(y | x)$ and a Transformer backbone that produces hidden states $z_t \in \mathbb{R}^D$ at each time step t for a given prompt x and generated response $y_{1:T}$.

E.1 Action vs. Semantic Exploration–Exploitation Metrics

In this subsection, we formalize the difference between the token-level *action* metrics used in prior RLHF/RLVR work and the hidden-state *semantic* metrics proposed in this paper. For a given prompt x and generated response $y_{1:T}$, let $\pi_\theta(\cdot | x, y_{<t})$ denote the model’s token-level policy distribution at step t , i.e. the softmax over the vocabulary induced by the logits at that position.

Token-level action space, log-probability, and entropy. We define the *average log probability* of an response and the *response entropy* as

$$\begin{aligned} \text{AvgLogProb}(x, y_{1:T}) &:= \frac{1}{T} \sum_{t=1}^T \log \pi_\theta(y_t | x, y_{<t}), \\ \text{RespEnt}(x, y_{1:T}) &:= \frac{1}{T} \sum_{t=1}^T H(\pi_\theta(\cdot | x, y_{<t})), \\ H(p) &:= - \sum_v p(v) \log p(v), \\ p(v) &:= \pi_\theta(v | x, y_{<t}). \end{aligned} \tag{34}$$

Thus *RespEnt* is the *token-level* entropy averaged over the response: at each step we compute the Shannon entropy of the vocabulary distribution and then average over time. At the response level semantic space, we consider x drawn from a prompt distribution \mathcal{P}_x and, for the purpose of analysis,

responses $y_{1:T}$ drawn *on-policy* from the model $p_\theta(\cdot | x)$:

$$\begin{aligned} \mathcal{L}_{\text{avg-log}}(\theta) &:= \\ &\mathbb{E}_{x \sim p(x)} \mathbb{E}_{y_{1:T} \sim p_\theta(\cdot | x)} \left[\text{AvgLogProb}(x, y_{1:T}) \right], \\ \mathcal{H}_{\text{resp}}(\theta) &:= \\ &\mathbb{E}_{x \sim p(x)} \mathbb{E}_{y_{1:T} \sim p_\theta(\cdot | x)} \left[\text{RespEnt}(x, y_{1:T}) \right]. \end{aligned} \tag{35}$$

Proposition E.1 (Token-level exploitation and exploration are tightly coupled). *Under on-policy sampling $y_{1:T} \sim p_\theta(\cdot | x)$, the corpus-level average log probability $\mathcal{L}_{\text{avg-log}}(\theta)$ and response entropy $\mathcal{H}_{\text{resp}}(\theta)$ satisfy*

$$\mathcal{L}_{\text{avg-log}}(\theta) = -\mathcal{H}_{\text{resp}}(\theta). \tag{36}$$

In particular, under the same sampling distribution, any change of the model that increases token-level exploitation in action space as measured by $\mathcal{L}_{\text{avg-log}}$ necessarily decreases $\mathcal{H}_{\text{resp}}$ by the same amount (and vice versa).

Proof. Fix a prompt x and a time step t . Conditioned on x and the history $y_{<t}$, the next token y_t is drawn from $\pi_\theta(\cdot | x, y_{<t})$. Taking the expectation of $\log \pi_\theta(y_t | x, y_{<t})$ under this distribution yields:

$$\begin{aligned} &\mathbb{E}_{y_t \sim \pi_\theta(\cdot | x, y_{<t})} \left[\log \pi_\theta(y_t | x, y_{<t}) \right] \\ &= \sum_v \pi_\theta(v | x, y_{<t}) \log \pi_\theta(v | x, y_{<t}) \\ &= -H(\pi_\theta(\cdot | x, y_{<t})). \end{aligned} \tag{37}$$

Now consider a full response $y_{1:T} \sim p_\theta(\cdot | x)$. By the law of iterated expectations,

$$\begin{aligned} &\mathbb{E}_{y_{1:T} \sim p_\theta(\cdot | x)} \left[\log \pi_\theta(y_t | x, y_{<t}) \right] \\ &= \mathbb{E}_{y_{<t} \sim p_\theta(\cdot | x)} \left[\mathbb{E}_{y_t \sim \pi_\theta(\cdot | x, y_{<t})} \left[\log \pi_\theta(y_t | x, y_{<t}) \right] \right] \\ &= -\mathbb{E}_{y_{<t} \sim p_\theta(\cdot | x)} \left[H(\pi_\theta(\cdot | x, y_{<t})) \right]. \end{aligned} \tag{38}$$

where we used Eq. 37 in the last step. Averaging

over $t = 1, \dots, T$ and dividing by T gives

$$\begin{aligned}
& \mathbb{E}_{y_{1:T} \sim p_\theta(\cdot|x)} \left[\text{AvgLogProb}(x, y_{1:T}) \right] \\
&= \mathbb{E}_{y_{1:T} \sim p_\theta(\cdot|x)} \left[\frac{1}{T} \sum_{t=1}^T \log \pi_\theta(y_t | x, y_{<t}) \right] \\
&= -\mathbb{E}_{y_{1:T} \sim p_\theta(\cdot|x)} \left[\frac{1}{T} \sum_{t=1}^T H(\pi_\theta(\cdot | x, y_{<t})) \right] \\
&= -\mathbb{E}_{y_{1:T} \sim p_\theta(\cdot|x)} \left[\text{RespEnt}(x, y_{1:T}) \right]. \quad (39)
\end{aligned}$$

Finally, taking expectation over prompts $x \sim p(x)$ on both sides of Eq. 39 yields

$$\mathcal{L}_{\text{avg-log}}(\theta) = -\mathcal{H}_{\text{resp}}(\theta), \quad (40)$$

which is exactly Eq. 36. This shows that under on-policy sampling, the two token-level metrics are related by a fixed negative sign and thus cannot be decoupled in the action space. \square

Hidden-state Effective Rank and velocity. The next proposition shows that these two hidden-state metrics in Sec. 3.1 and 3.3 are *structurally decoupled* at the level of trajectories: knowing the final Effective Rank alone does not determine ERV, and conversely.

Proposition E.2 (Hidden-state metrics are structurally decoupled). *Fix $K \geq 3$ evaluation steps. Consider the map that associates to each Effective Rank trajectory $m = (m_1, \dots, m_K) \in \mathbb{R}^K$ its final value*

$$\text{ER}_{\text{final}}(m) := m_K \quad (41)$$

and its Effective Rank velocity

$$\text{ERV}(m) := \frac{1}{K-1} \sum_{j=2}^K \left(m_j - \frac{1}{j-1} \sum_{k=1}^{j-1} m_k \right). \quad (42)$$

Then:

1. *There is no function $f : \mathbb{R} \rightarrow \mathbb{R}$ such that $\text{ERV}(m) = f(\text{ER}_{\text{final}}(m))$ for all trajectories $m \in \mathbb{R}^K$.*
2. *There is no function $g : \mathbb{R} \rightarrow \mathbb{R}$ such that $\text{ER}_{\text{final}}(m) = g(\text{ERV}(m))$ for all trajectories $m \in \mathbb{R}^K$.*

Equivalently, ER_{final} and ERV are not functionally dependent: they capture genuinely different aspects of the Effective Rank sequence.

Proof. We view ER_{final} and ERV as real-valued functions on \mathbb{R}^K . The proof is purely algebraic and does not rely on any monotonicity of m_j .

Step 1: ERV is a non-trivial linear functional. Introduce the shorthand

$$\Delta m_j := m_j - m_{j-1}, \quad j \geq 2. \quad (43)$$

A direct calculation shows that each increment δ_j can be written as

$$\delta_j = \frac{1}{j-1} \sum_{r=2}^j (r-1) \Delta m_r, \quad j \geq 2, \quad (44)$$

so that ERV is a linear functional of m :

$$\text{ERV}(m) = \sum_{j=1}^K \alpha_j m_j, \quad (45)$$

for some fixed coefficients $\alpha_1, \dots, \alpha_K$ that depend only on K (and s) and satisfy $\sum_{j=1}^K \alpha_j = 0$ and $\alpha_j \neq 0$ for at least two indices j (e.g. $\alpha_1 \neq 0$ and $\alpha_K \neq 0$). In particular, ERV is *not* proportional to the projection onto any single coordinate m_j .

Step 2: No functional dependence of ERV on the final ER. Suppose, for contradiction, that there exists a function $f : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$\begin{aligned}
\text{ERV}(m) &= f(\text{ER}_{\text{final}}(m)) \\
&= f(m_K) \quad \text{for all } m \in \mathbb{R}^K.
\end{aligned} \quad (46)$$

Fix any constant $c \in \mathbb{R}$. Consider the affine subspace

$$\mathcal{A}_c := \{m \in \mathbb{R}^K : m_K = c\}. \quad (47)$$

On this subspace, $\text{ER}_{\text{final}}(m) \equiv c$ is constant, so by assumption $\text{ERV}(m) \equiv f(c)$ must also be constant. However, ERV is a non-trivial linear functional that depends on at least one coordinate m_j with $j < K$. Therefore, restricted to \mathcal{A}_c , the map $m \mapsto \text{ERV}(m)$ varies with those coordinates and cannot be constant. This yields a contradiction. Hence no such f exists.

Step 3: No functional dependence of final ER on ERV. The argument is symmetric. Suppose there exists $g : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$\begin{aligned}
m_K &= \text{ER}_{\text{final}}(m) \\
&= g(\text{ERV}(m)) \quad \text{for all } m \in \mathbb{R}^K.
\end{aligned} \quad (48)$$

Fix any constant $c \in \mathbb{R}$ and consider the affine subspace

$$\mathcal{B}_c := \{m \in \mathbb{R}^K : \text{ERV}(m) = c\}. \quad (49)$$

Since ERV is a non-trivial linear functional, \mathcal{B}_c is an affine hyperplane of codimension 1, and m_K can vary freely among its points. Yet the assumed relation $m_K = g(\text{ERV}(m)) = g(c)$ would force m_K to be constant on \mathcal{B}_c , which is impossible. Thus no such g exists.

Combining the two steps, we conclude that ER_{final} and ERV are not functionally dependent on each other. \square

Summary. Proposition E.1 establishes that two classical token-level metrics, average log probability and response entropy, are *algebraically coupled* under on-policy sampling and therefore cannot be varied independently. In contrast, Proposition E.2 shows that our hidden-state metrics, terminal Effective Rank and Effective Rank Velocity, are structurally decoupled. They depend on different functionals of the Effective Rank trajectory and capture complementary aspects of exploration (semantic diversity level) and exploitation (the rate of semantic diversity gain).

E.2 Effective Rank as Semantic Exploration

We now formalize the interpretation of Effective Rank as a measure of semantic diversity and uncertainty in the hidden-state space, and hence as a representation-level proxy for exploration in LLM reasoning. We assume that the hidden states are semantic representations in the sense that downstream semantic properties can be approximately recovered as linear functionals of the hidden vectors. This is standard in representation learning and probing work on large language models.

Assumption E.3 (Semantic linear decodability). There exists a collection of K semantic features $s^{(1)}, \dots, s^{(K)}$ (e.g., semantic roles, entity identities, factual attributes, intermediate reasoning states) such that for each time step t and feature index k we have

$$s_t^{(k)} \approx w_k^\top h_t, \quad w_k \in \mathbb{R}^D. \quad (50)$$

Semantic features are therefore approximately linearly decodable from hidden states.

Assumption E.4 (Bounded energy and finite support). For a given trajectory $Z_{1:T}$, there exists an orthonormal basis of semantic directions $\{e_1, \dots, e_D\}$ such that each hidden state admits a decomposition

$$z_t = \sum_{i=1}^D a_{t,i} e_i, \quad (51)$$

with $\sum_{t=1}^T a_{t,i}^2 < \infty$ for all i , and only finitely many coordinates $a_{t,i}$ carry task-relevant semantic variation.

Proposition E.5 (Effective Rank as semantic diversity and uncertainty). *Let $Z_{1:T}$ have singular values $\{\sigma_i\}_{i=1}^{\min(T,D)}$ (i.e., the diagonal of Σ in its SVD), and let $\text{ER}(Z_{1:T})$ denote its Effective Rank.*

1. *If the trajectory uses exactly k orthogonal semantic directions with equal energy and no others, i.e. the singular values satisfy $\sigma_1 = \dots = \sigma_k > 0$ and $\sigma_{k+1} = \dots = \sigma_r = 0$, then $\text{ER}(Z_{1:T}) = k$.*
2. *More generally, if the singular value spectrum of $Z_{1:T}$ becomes more spread out over more directions in the sense of majorization (i.e. the normalized singular value vector becomes more uniform over a larger support), then $\text{ER}(Z_{1:T})$ increases.*

Consequently, $\text{ER}(Z_{1:T})$ is a basis-invariant, strictly increasing measure of the number of independent semantic directions that are effectively used by the hidden states, and thus a natural representation-level proxy for semantic exploration and uncertainty.

Proof. We proceed in two parts.

(1) Equal-energy k -dimensional semantic subspace. Suppose $Z_{1:T}$ uses exactly k orthogonal semantic directions with equal energy. Then, up to permutation, the non-zero singular values satisfy

$$\sigma_1 = \dots = \sigma_k = c > 0, \quad \sigma_{k+1} = \dots = \sigma_r = 0. \quad (52)$$

The normalized singular values are thus

$$q_i = \begin{cases} 1/k, & i = 1, \dots, k, \\ 0, & i > k, \end{cases} \quad (53)$$

and the entropy of q is

$$H(q) = - \sum_{i=1}^k \frac{1}{k} \log \frac{1}{k} = \log k. \quad (54)$$

Therefore

$$\text{ER}(Z_{1:T}) = \exp(H(q)) = \exp(\log k) = k. \quad (55)$$

This shows that, in the idealized case of exactly k equi-energetic semantic directions, Effective Rank matches the true semantic dimensionality k .

(2) Monotonicity under majorization. Consider two hidden-state trajectories Z and \tilde{Z} with normalized singular value spectra q and \tilde{q} , respectively. Suppose that q is majorized by \tilde{q} (denoted $q < \tilde{q}$), meaning intuitively that \tilde{q} is “more spread out” and therefore more uniform across a larger support.

It is a standard result in information theory that the Shannon entropy $H(\cdot)$ is Schur-concave: if $q < \tilde{q}$, then $H(q) \leq H(\tilde{q})$ with strict inequality whenever $q \neq \tilde{q}$. Therefore

$$\text{ER}(Z) = \exp(H(q)) \leq \exp(H(\tilde{q})) = \text{ER}(\tilde{Z}), \quad (56)$$

with strict inequality when the majorization is strict. In words, whenever the singular value spectrum becomes more spread out across more directions, the Effective Rank strictly increases.

Summary. Combining the two parts, we see that Effective Rank equals the number of equi-energetic semantic directions in the idealized case and increases whenever the representation distributes energy over more orthogonal directions. Since, by Assumption E.3, semantic features are linearly decodable along such directions, $\text{ER}(Z_{1:T})$ provides a basis-invariant measure of how many independent semantic dimensions are explored by the hidden states and how evenly they are used. This justifies its interpretation as a representation-level exploration and uncertainty metric. \square

E.3 Effective Rank Velocity as Semantic Exploitation

Building on Sec. 3.1, where Effective Rank (ER) is shown to measure the number and uniform use of semantic directions in the hidden-state space, we now give a representation-only justification for interpreting Effective Rank Velocity (ERV) as *semantic exploitation*.

Throughout this subsection we fix a single trajectory $Z_{1:T}$ and a stride s . Let the evaluation positions be $t_j = js$ for $j = 1, \dots, K$ with $K = \lfloor (T-1)/s \rfloor$, and write

$$m_j := \text{ER}(Z_{1:t_j}), \quad j = 1, \dots, K. \quad (57)$$

Thus, $\{m_j\}_{j=1}^K$ is the ER trajectory of the growing prefixes of the same response.

ERV as a recency-weighted sum of ER increments. For convenience we recall the notation of Def. 3.3 with $M = \text{ER}$. Define the local ER increments

$$\Delta m_r := m_r - m_{r-1}, \quad r \geq 2. \quad (58)$$

Def. 3.3 introduces the “instantaneous difference”

$$\delta_j := m_j - \frac{1}{j-1} \sum_{k=1}^{j-1} m_k, \quad j \geq 2, \quad (59)$$

and the first-order temporal difference (ERV) as

$$\Delta_{\text{ER}}^{(1)} := \frac{1}{K-1} \sum_{j=2}^K \delta_j. \quad (60)$$

The following lemma makes explicit that ERV is a recency-weighted average of the consecutive ER increments.

Lemma E.6 (Recency-weighted velocity of ER). *For any sequence $(m_j)_{j=1}^K$, the instantaneous differences admit the representation*

$$\delta_j = \frac{1}{j-1} \sum_{r=2}^j (r-1) \Delta m_r, \quad j \geq 2, \quad (61)$$

and hence ERV can be written as

$$\begin{aligned} \Delta_{\text{ER}}^{(1)} &= \sum_{r=2}^K w_r \Delta m_r, \\ w_r &:= \frac{r-1}{K-1} \sum_{j=r}^K \frac{1}{j-1} > 0. \end{aligned} \quad (62)$$

In particular, ERV is a positive linear combination of the local ER increments Δm_r , assigning larger weights to more recent steps.

Proof. Eq. 61 is exactly in Def. 3.3 with $M = \text{ER}$, obtained by expressing δ_j in terms of the increments Δm_r via telescoping. Plugging Eq. 61 into the definition of $\Delta_{\text{ER}}^{(1)}$ and exchanging the order of summation yields

$$\begin{aligned} \Delta_{\text{ER}}^{(1)} &= \frac{1}{K-1} \sum_{j=2}^K \frac{1}{j-1} \sum_{r=2}^j (r-1) \Delta m_r \\ &= \sum_{r=2}^K \left[\frac{r-1}{K-1} \sum_{j=r}^K \frac{1}{j-1} \right] \Delta m_r. \end{aligned} \quad (63)$$

which gives Eq. 62 with w_r as stated. Since $r-1 > 0$ and the harmonic tail $\sum_{j=r}^K (j-1)^{-1}$ is positive, we have $w_r > 0$ for all r . \square

Semantic exploitation as positive ER drift in a fixed semantic subspace.

Sec. E.1 states that, up to an orthonormal change of basis, the hidden states can be written as $h_t = \sum_i a_{t,i} e_i$ with bounded energy along each semantic direction e_i , and that ER is a strictly increasing, basis-invariant measure of how many semantic directions are effectively used and how evenly energy is distributed among them. In particular, if the set of active directions (support of the singular value spectrum) is kept fixed and the spectrum becomes more uniform (in the sense of majorization), then ER strictly increases. Motivated by this, we isolate an idealized *semantic exploitation* regime in which the trajectory has already selected a semantic subspace and is refining it.

Definition E.7 (Semantic exploitation regime). Let $(m_j)_{j=1}^K$ be the ER trajectory of a response, and let $q^{(j)}$ denote the normalized singular value vector of $Z_{1:t_j}$. We say that steps $j = 2, \dots, K$ form a *semantic exploitation regime with rate $\mu > 0$* if:

1. (Fixed semantic support) The support of $q^{(j)}$ is independent of j , i.e., the set of active semantic directions is fixed.
2. (Uniformization within the support) For every $j \geq 2$, $q^{(j)}$ is more uniform than $q^{(j-1)}$ on this fixed support, in the sense of majorization, so that by Prop. E.5 we have $m_j - m_{j-1} = \Delta m_j \geq \mu$ for some $\mu > 0$.

Intuitively, condition (i) says that the model has committed to a particular semantic subspace (a line of reasoning), and condition (ii) says that it keeps redistributing energy within this subspace to make use of all its semantic directions more evenly. This is precisely the notion of “refining a promising strategy” in representation space.

ERV lower-bounds the semantic exploitation rate. Under Def. E.7, ER experiences a persistent positive drift along the trajectory. The next proposition shows that ERV is a quantitative lower bound on this drift, and thus a natural measure of semantic exploitation strength.

Proposition E.8 (ERV as a lower bound on semantic exploitation rate). *Assume the hidden states satisfy Assumptions E.3 and E.4 and that steps $j = 2, \dots, K$ form a semantic exploitation regime with rate $\mu > 0$ in the sense of Def. E.7, so that $\Delta m_j \geq \mu$ for all $j \geq 2$. Then*

$$\Delta_{\text{ER}}^{(1)} \geq \frac{\mu K}{4}. \quad (64)$$

In particular, ERV is strictly positive and grows linearly with the length K of the exploitation segment.

Proof. By Eq. 61 and the assumption $\Delta m_r \geq \mu$ we obtain, for each $j \geq 2$,

$$\begin{aligned} \delta_j &= \frac{1}{j-1} \sum_{r=2}^j (r-1) \Delta m_r \\ &\geq \frac{1}{j-1} \sum_{r=2}^j (r-1) \mu \\ &= \frac{\mu}{j-1} \sum_{r=2}^j (r-1) = \frac{\mu j}{2}. \end{aligned} \quad (65)$$

Averaging over j then yields

$$\begin{aligned} \Delta_{\text{ER}}^{(1)} &= \frac{1}{K-1} \sum_{j=2}^K \delta_j \\ &\geq \frac{1}{K-1} \sum_{j=2}^K \frac{\mu j}{2} = \frac{\mu}{2(K-1)} \sum_{j=2}^K j. \end{aligned} \quad (66)$$

For $K \geq 2$ we have $\sum_{j=2}^K j \geq \frac{K(K-1)}{2}$, so

$$\Delta_{\text{ER}}^{(1)} \geq \frac{\mu}{2(K-1)} \cdot \frac{K(K-1)}{2} = \frac{\mu K}{4}, \quad (67)$$

which proves Eq. 64. \square

Thus, in an idealized regime where the model has already discovered a useful semantic subspace and is consistently enriching it, ERV provides a strictly positive, linearly growing lower bound on the rate at which semantic complexity within that subspace is being exploited.

F Implementation Details

F.1 Training and Evaluation Details

We typically use the same set of hyperparameters to train and evaluate all models in the SimpleRL-Zoo series (Zeng et al., 2025) in the default main experiment setting.

Training. We conduct all experiments with 4 A800-PCIE-80GB GPUs. For GRPO and PPO, we use a prompt batch size of 48 with a maximum rollout length of 1536 tokens. Training is performed using a mini-batch size of 24. For GRPO, we generate 4 rollouts per prompt. For PPO, we use DeepSeek-R1-Distill-Qwen-1.5B (Guo et al., 2025) as the value model and generate 1 rollout per prompt. The default sampling temperature is set to 1.0, and the clip ratio is 0.2. For all actor

models ranging from 3B to 8B parameters, we use a learning rate of 1e-6 and a KL loss coefficient of 1e-4. For critic models in PPO, we use a learning rate of 1e-5. For our training datasets, we follow the same setup as in Zeng et al. (2025), where the data is filtered from GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021b) configured with different difficulty levels for models of varying capabilities. We tested using the checkpoint model trained up to 120 steps.

Evaluation. We build our evaluation script based on that of Zeng et al. (2025), using a temperature of 0.6 and a maximum generation length of 2048 tokens. To ensure consistency, we adopt the same prompt template used during training. For most benchmarks, we report Pass@1 results. However, for benchmarks like AIME 2024, which contains fewer problems, we report both Pass@1 and average accuracy (Pass@256), computed over 256 generated samples per problem.

Base Models. To demonstrate the universality of our insights and methods, we conduct zero RL training experiments on Llama-3.2 (3B), Llama-3.1 (8B) (Dubey et al., 2024), Mistral-v0.3-7B (Jiang et al., 2024), and Qwen-2.5 (1.5B, 3B, 7B) (Hui et al., 2024). For value model in PPO, we use DeepSeek-R1-Distill-Qwen-1.5B (Guo et al., 2025) for all experiments.

Benchmark. We evaluate on a diverse suite of mathematical reasoning benchmarks. These include standard benchmarks such as GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021b), ASDiv (Miao et al., 2020), Carp (English Version) (Zhang et al., 2023), MAWPS (Koncel-Kedziorski et al., 2016), SVAMP (Patel et al., 2021), TabMWP (Lu et al., 2023), and OlympiadBench (He et al., 2024); Chinese mathematics collections like CMATH (Wei et al., 2023) and Gaokao 2024; and benchmarks from mathematics competitions, including the 2024/2025 AIME and the 2023/2024 AMC.

F.2 Efficient Incremental Computation of Higher-Order Metrics

A naive computation of the temporal difference metrics would be computationally prohibitive. Our method’s feasibility hinges on an efficient, incremental algorithm that computes the required metrics without redundant operations on the growing hidden state matrix $\mathbf{Z} \in \mathbb{R}^{T \times D}$.

The effective rank is derived from the singu-

lar values of the mean-centered hidden state matrix \mathbf{Z}_c . These are equivalent to the square roots of the eigenvalues of the centered Gram matrix $\mathbf{K} = \mathbf{Z}_c \mathbf{Z}_c^\top$. Instead of recomputing \mathbf{K}_t from scratch at each time step t , our algorithm incrementally updates two sufficient statistics: the uncentered Gram matrix $\mathbf{U}_t = \mathbf{Z}_{1:t} \mathbf{Z}_{1:t}^\top$ and the sum of hidden state vectors $s_t = \sum_{i=1}^t z_i$. When extending the analysis window, the new uncentered Gram matrix \mathbf{U}_t is constructed from the prior matrix \mathbf{U}_{t-s} and the new chunk of hidden states $\Delta \mathbf{Z}_t = \mathbf{Z}_{t-s+1:t}$. This update follows a recursive block matrix structure:

$$\mathbf{U}_t = \begin{pmatrix} \mathbf{U}_{t-s} & \mathbf{Z}_{1:t-s} (\Delta \mathbf{Z}_t)^\top \\ (\Delta \mathbf{Z}_t) \mathbf{Z}_{1:t-s}^\top & (\Delta \mathbf{Z}_t) (\Delta \mathbf{Z}_t)^\top \end{pmatrix} \quad (68)$$

From the efficiently updated \mathbf{U}_t and s_t , we can directly construct the centered Gram matrix \mathbf{K}_t . Letting $\boldsymbol{\mu}_t = s_t/t$ be the mean vector and $\mathbf{1}_t$ be a column vector of ones, \mathbf{K}_t can be expressed as:

$$\mathbf{K}_t = \mathbf{U}_t - (\mathbf{Z}_{1:t} \boldsymbol{\mu}_t) \mathbf{1}_t^\top - \mathbf{1}_t (\mathbf{Z}_{1:t} \boldsymbol{\mu}_t)^\top + (\boldsymbol{\mu}_t^\top \boldsymbol{\mu}_t) \cdot (\mathbf{1}_t \mathbf{1}_t^\top). \quad (69)$$

This allows for the calculation of \mathbf{K}_t without re-accessing the full history of hidden states. The eigenvalues $\{\lambda_j\}$ of \mathbf{K}_t are then used to derive the effective rank. First, the singular values of the centered matrix are obtained, $\sigma_j = \sqrt{\lambda_j}$. These are normalized to form a probability distribution, $p_j = \sigma_j / \sum_k \sigma_k$. The effective rank is then the exponential of the Shannon entropy of this distribution: $\text{erank}(\mathbf{Z}_{c,t}) = \exp\left(-\sum_j p_j \log p_j\right)$. This pipeline efficiently yields a sequence of effective rank values, $m_{j \cdot s} = \text{erank}(\mathbf{Z}_{c,j \cdot s})$, at each stride s . From this sequence, we compute the instantaneous first-order difference δ , which compares the current value to the running average of all preceding values. This is defined recursively as: $\delta_{j \cdot s} = m_{j \cdot s} - \frac{1}{j-1} \sum_{k=1}^{j-1} m_{k \cdot s}$.

The computational advantage of this incremental approach is substantial. While the total cost for the series of eigenvalue decompositions $\mathcal{O}(T^4/s)$, is common to both methods, the cost of matrix construction differs significantly. The naive method’s recalculation totals $\mathcal{O}(DT^3/s)$, whereas our incremental update method reduces this to $O(DT^2)$. This reduction of the polynomial dependency on sequence length T from cubic to quadratic is critical, as this term is multiplied by the large hidden dimension D , making it the dominant factor

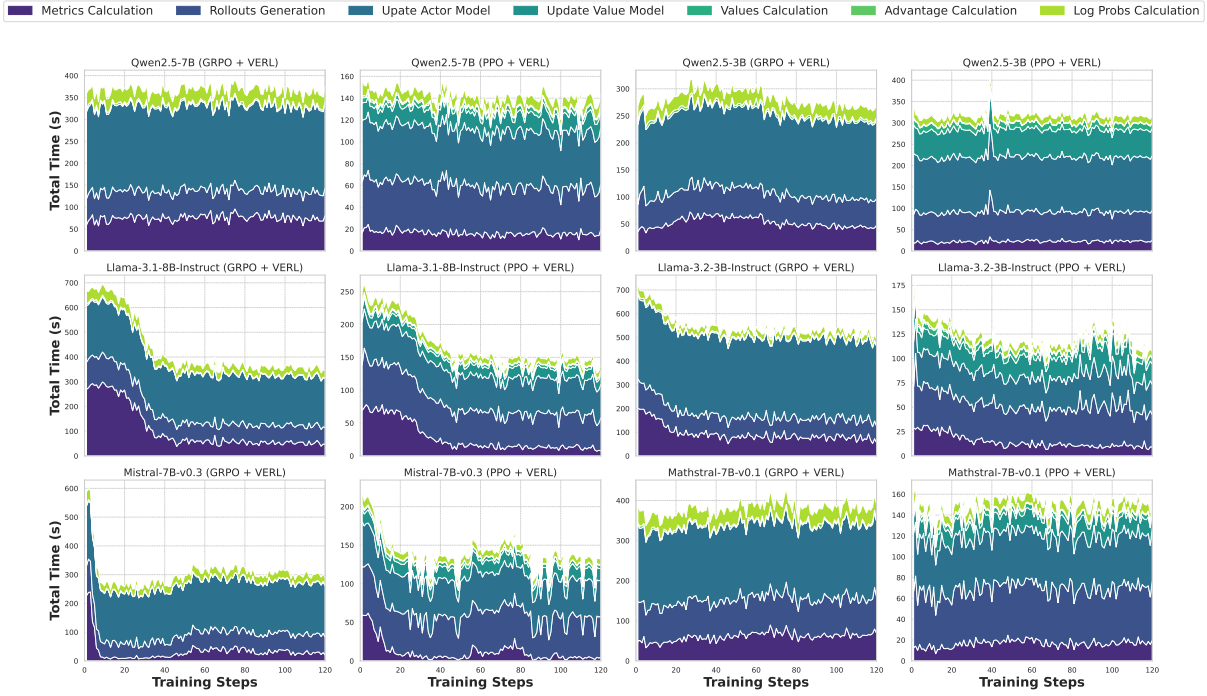


Figure 6: Time overhead of the main computation of RL Training.

in practical performance and rendering the dense calculation of temporal dynamics feasible. In the *worst-case* scenario where the sequence length T exceeds the hidden dimension D , and both D and the stride s can be treated as constants. The naïve approach that reconstructs matrices independently at each stride has a matrix-construction cost scaling as $\mathcal{O}(T^2)$, VERL’s incremental Gram/covariance updates scale as $\mathcal{O}(T)$. So asymptotically, our implementation has a strictly better dependency on T than a naïve SVD-based design.

F.3 Time Overhead of VERL Training

We conducted post-training with Zero RL on several base models. The Fig. 6 illustrates the time associated with each computational stage. The ‘metrics calculation’ component, which represents the cost of computing metrics for hidden states, accounts for an insignificant portion of the total processing time. This demonstrates that our method does not introduce substantial time overhead. To further stress-test the *worst-case* scenario, we deliberately compute ER, ERV, and ERA on the CPU rather than the GPU, and still observe that the additional time overhead remains negligible.

G More Experiments

G.1 Detailed Analysis of Response-Level Hidden States

As shown in Figs. 7 and 8, our analysis of response-level hidden states across additional LLMs con-

firms that the insights presented in Sec. 4.1 hold true for various base models and RL paradigms.

G.2 Detailed Analysis of Dataset-Level Hidden States

As shown in Figs. 9 and 10, our analysis of dataset-level hidden states across additional LLMs confirms that the insights presented in Sec. 4.2 hold for various base models and RL paradigms.

G.3 Detailed Analysis of Pass@1 Performance

As shown in Tab. 4, Pass@1 measures the model’s ability to generate a correct answer in a single attempt, which directly reflects its exploitation ability. We fine-tune the base model by integrating our VERL-based Advantage method into two reinforcement learning paradigms, GRPO and PPO.

G.4 Detailed Analysis of Pass@ k Performance

As shown in Tab. 5, which provides a comprehensive analysis of the models’ performance on the Pass@ k metric, which is a direct measure of the model’s exploration ability. As a supplement to the main paper’s discussion, it presents the detailed performance of our VERL-based models across a variety of mathematical reasoning benchmarks. These results demonstrate that our method consistently leads to significant improvements, confirming its effectiveness in enhancing the models’ exploration capabilities.

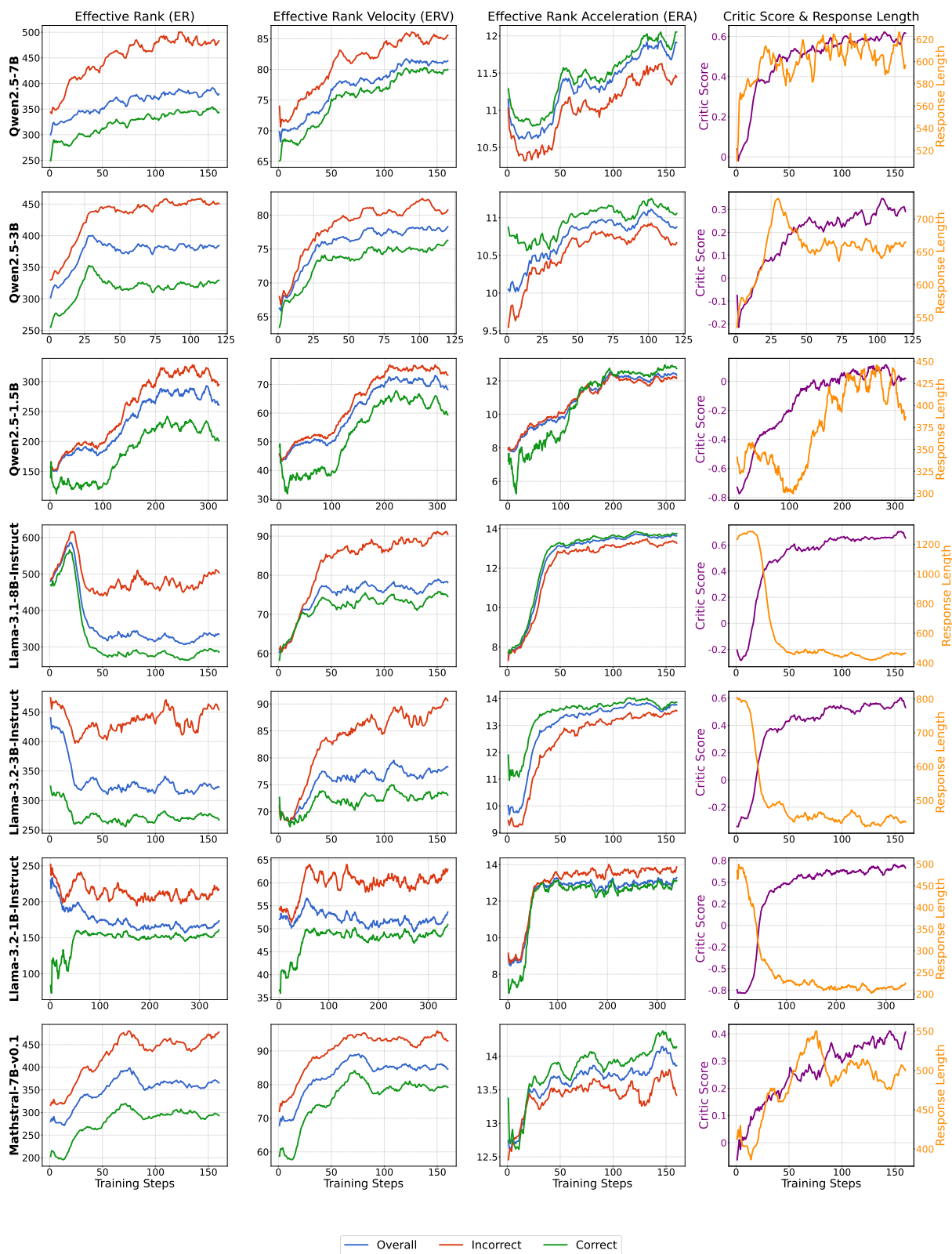


Figure 7: Response-level metrics during GRPO post-training, smoothed with a 10-step rolling window. Metrics are shown for the Overall batch, as well as for subsets of Correct and Incorrect samples. The rightmost column displays the average Critic Score (reward) and Response Length per batch.

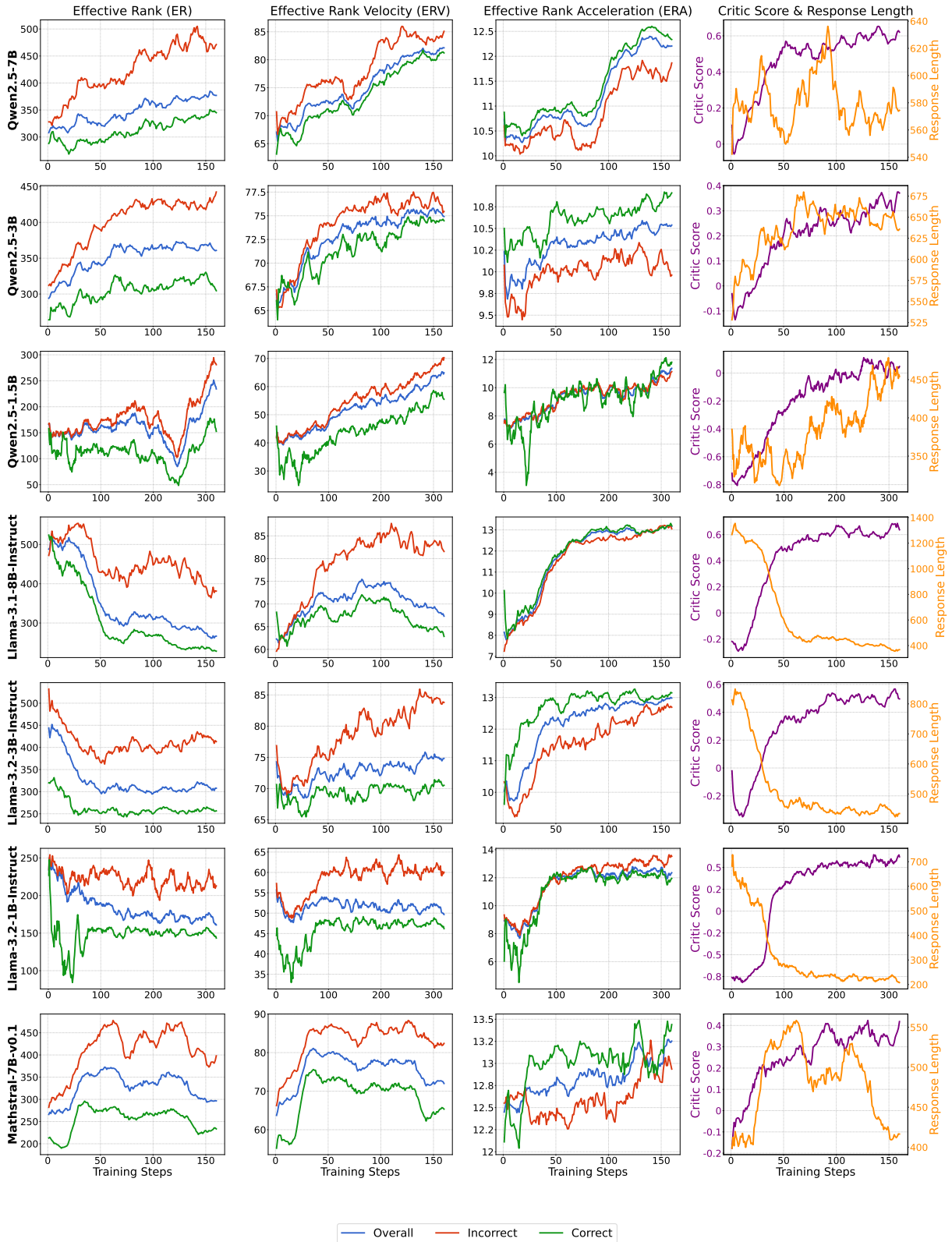


Figure 8: Response-level metrics during PPO post-training, smoothed with a 10-step rolling window. Metrics are shown for the Overall batch, as well as for subsets of Correct and Incorrect samples. The rightmost column displays the average Critic Score (reward) and Response Length per batch.

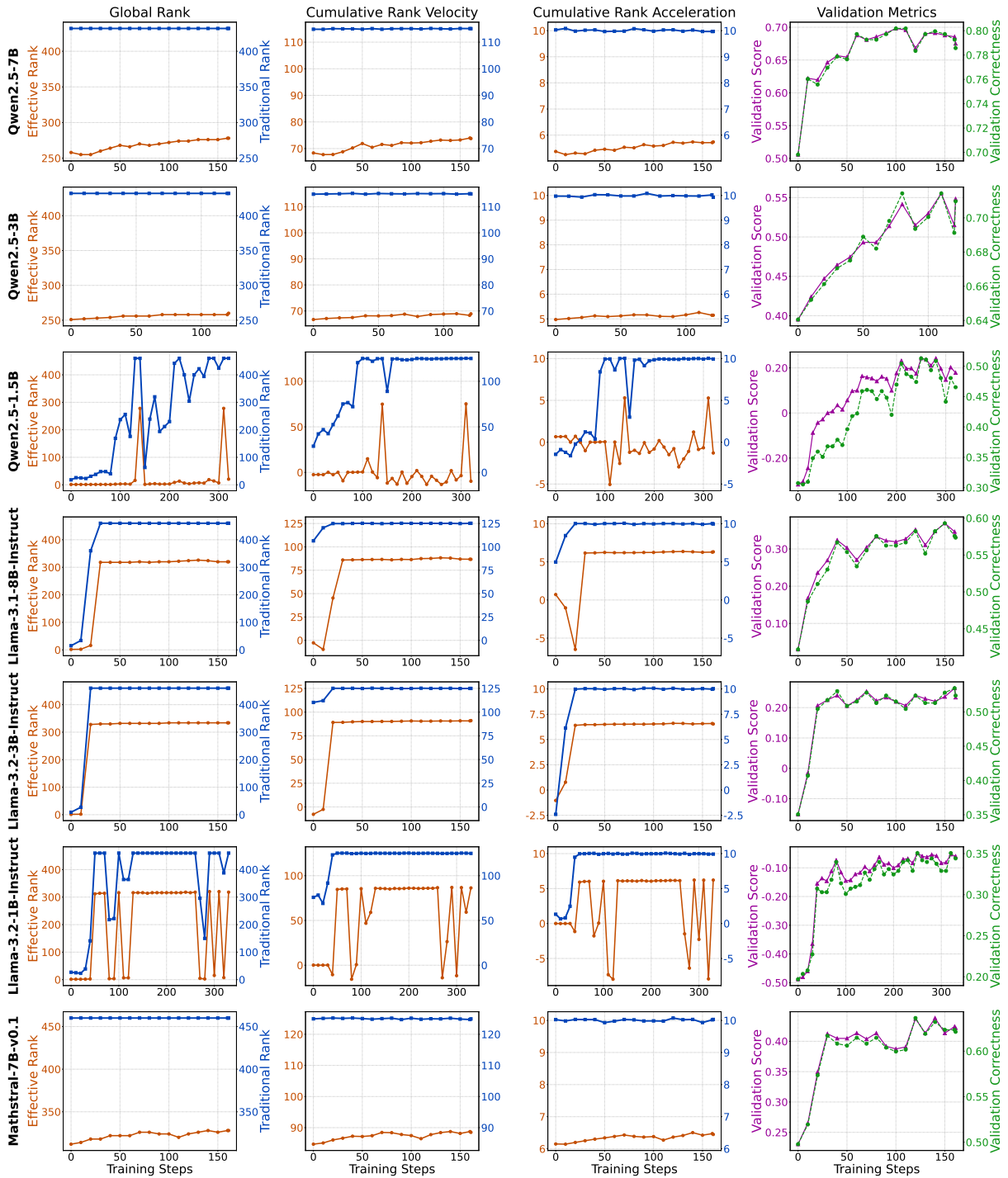


Figure 9: Visualization of dataset-level metrics during GRPO post-training. The figure compares **Traditional** metrics with our **proposed** metrics. Also shown are the **Validation Score** and sample **Correctness**, both averaged over the validation dataset.

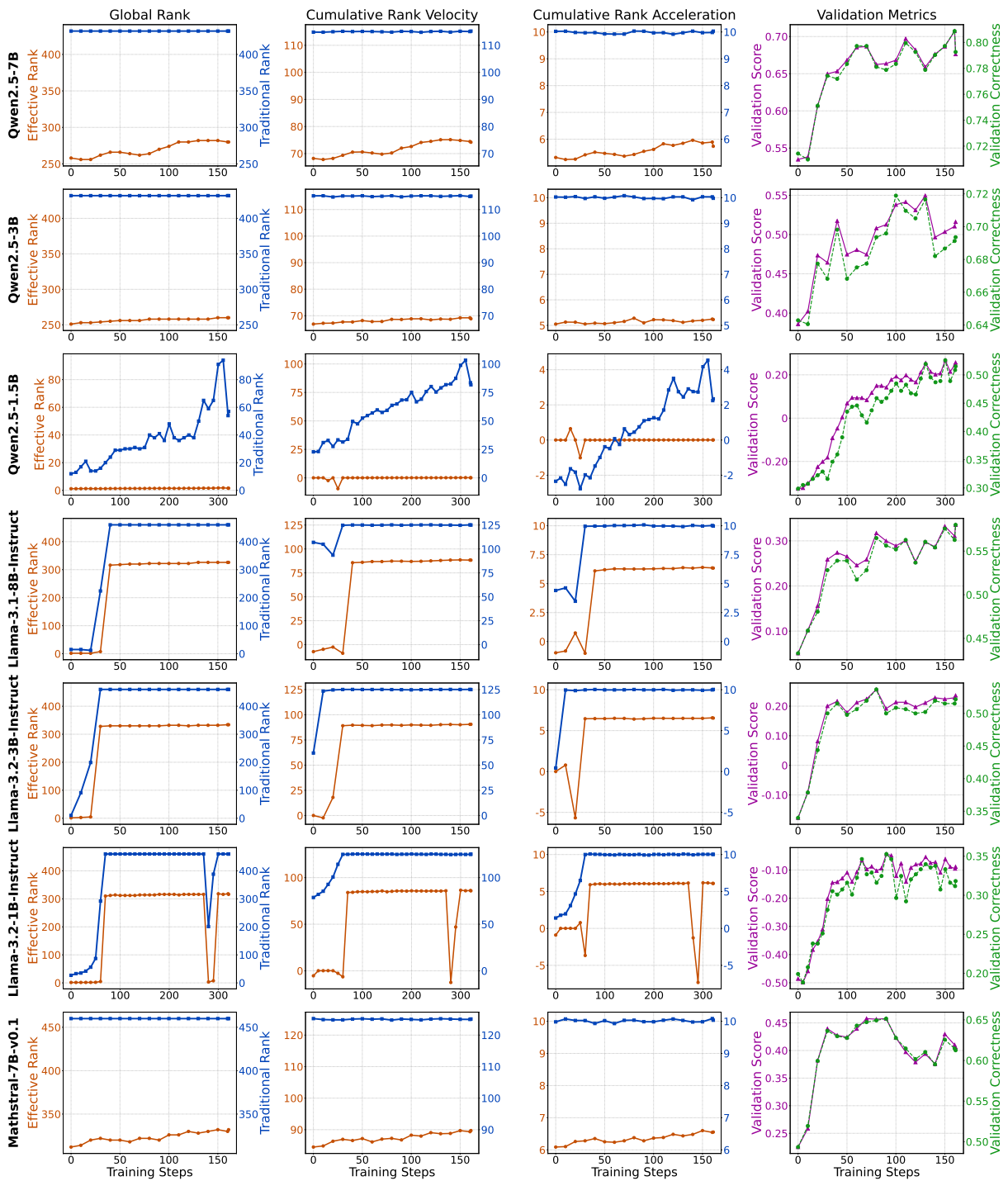


Figure 10: Visualization of dataset-level metrics during PPO post-training. The figure compares **Traditional** metrics with our **proposed** metrics. Also shown are the **Validation Score** and sample **Correctness**, both averaged over the validation dataset.

Table 4: Performance comparison of instruction-tuned models on mathematical reasoning benchmarks (Pass@1). “+ GRPO” and “+ PPO” denote reinforcement learning fine-tuning from the base model using GRPO and PPO, respectively. “w/ VERL.” indicates the application of our VERL-based advantage to the corresponding RL algorithm. Δ represents the performance difference between the baseline RL method and its VERL-advanced variant. All results are reported in percentage (%).

Model	AIME24	AIME25	AMC23	AMC24	ASDiv	Carp_En	CMATH	Gaokao 2024_I	Gaokao 2024_Mix	Gaokao MathCloze	GSM8K	MAWPS	Olympiad Bench	SVAMP	TabMWP	Avg.
Llama-3.2-3B-Instruct	0.0	0.0	25.0	11.1	74.6	26.5	10.2	14.3	14.3	6.8	66.6	86.9	12.7	74.1	41.4	31.0
+ GRPO	3.3	0.0	27.5	8.9	88.8	45.0	28.3	21.4	20.9	23.7	80.7	96.0	16.7	87.7	71.7	41.4
+ GRPO w/ VERL.	13.3	6.7	25.0	11.1	89.3	45.4	46.2	14.3	22.0	22.9	81.7	96.0	17.6	87.8	72.3	43.4
Δ_{GRPO}	+10.0	+6.7	-2.5	+2.2	+0.5	+0.4	+17.9	-7.1	+1.1	-0.8	+1.0	+0.0	+0.9	+0.1	+0.6	+2.0
+ PPO	10.0	3.3	22.5	13.3	87.9	46.4	21.2	7.1	16.5	20.3	81.4	95.5	17.8	86.8	71.0	40.1
+ PPO w/ VERL.	10.0	3.3	25.0	11.1	88.7	46.0	30.7	14.3	19.8	27.1	82.9	95.7	17.3	85.8	71.3	41.9
Δ_{PPO}	+0.0	+0.0	+2.5	-2.2	+0.8	-0.4	+9.5	+7.2	+3.3	+6.8	+1.5	+0.2	-0.5	-1.0	+0.3	+1.9
Llama-3.1-8B-Instruct	0.0	3.3	17.5	8.9	48.0	34.1	18.5	0.0	15.4	16.9	47.4	43.5	10.4	48.5	34.3	23.1
+ GRPO	6.7	0.0	22.5	15.6	90.3	42.4	60.7	7.1	14.3	32.2	88.4	96.4	19.7	88.5	82.7	44.5
+ GRPO w/ VERL.	10.0	3.3	32.5	15.6	90.7	45.0	72.7	14.3	14.3	30.5	88.6	96.9	21.3	88.4	83.1	47.2
Δ_{GRPO}	+3.3	+3.3	+10.0	+0.0	+0.4	+2.6	+12.0	+7.2	+0.0	-1.7	+0.2	+0.5	+1.6	-0.1	+0.4	+2.7
+ PPO	6.7	0.0	30.0	17.8	89.8	42.0	60.0	0.0	14.3	25.4	86.4	95.7	18.2	88.6	82.3	43.8
+ PPO w/ VERL.	10.0	0.0	35.0	13.3	90.7	42.6	62.0	14.3	22.0	28.8	87.3	96.6	19.1	88.1	83.0	46.2
Δ_{PPO}	+3.3	+0.0	+5.0	-4.5	+0.9	+0.6	+2.0	+14.3	+7.7	+3.4	+0.9	+0.9	+0.9	-0.5	+0.7	+2.4
Qwen2.5-3B	6.7	0.0	20.0	24.4	90.7	54.7	76.7	0.0	22.0	41.5	80.7	95.1	23.0	84.3	71.3	46.1
+ GRPO	3.3	0.0	40.0	22.2	92.6	56.0	82.7	7.1	27.5	42.4	82.8	96.6	23.6	89.0	81.4	49.8
+ GRPO w/ VERL.	6.7	0.0	30.0	17.8	92.6	56.9	84.8	21.4	33.0	49.2	82.2	96.4	24.4	88.5	81.0	51.0
Δ_{GRPO}	+3.4	+0.0	-10.0	-4.4	+0.0	+0.9	+2.1	+14.3	+5.5	+6.8	-0.6	-0.2	+0.8	-0.5	-0.4	+1.2
+ PPO	3.3	0.0	32.5	15.6	92.8	56.5	83.2	0.0	28.6	50.0	81.7	96.6	24.4	86.0	80.8	48.8
+ PPO w/ VERL.	6.7	0.0	32.5	17.8	92.6	57.0	84.3	21.4	29.7	47.5	81.8	96.5	24.6	88.3	81.4	50.8
Δ_{PPO}	+3.4	+0.0	+0.0	+2.2	-0.2	+0.5	+1.1	+21.4	+1.1	-2.5	+0.1	-0.1	+0.2	+2.3	+0.6	+2.0
Qwen2.5-7B	6.7	0.0	45.0	15.6	91.4	55.8	86.7	42.9	33.0	49.2	85.8	95.4	25.8	88.5	82.8	53.6
+ GRPO	10.0	6.7	55.0	26.7	94.8	60.2	91.7	14.3	34.1	64.4	90.2	97.6	36.1	92.8	91.3	57.7
+ GRPO w/ VERL.	13.3	10.0	50.0	28.9	95.0	60.8	90.7	35.7	35.2	69.5	89.2	97.7	35.4	92.9	91.9	59.8
Δ_{GRPO}	+3.3	+3.3	-5.0	+2.2	+0.2	+0.6	-1.0	+21.4	+1.1	+5.1	-1.0	+0.1	-0.7	+0.1	+0.6	+2.1
+ PPO	6.7	3.3	50.0	33.3	94.9	59.6	89.8	28.6	31.9	63.6	89.1	97.3	36.1	92.8	90.8	57.9
+ PPO w/ VERL.	10.0	6.7	52.5	33.3	94.8	60.0	90.3	28.6	34.1	66.9	90.2	97.8	36.1	92.5	90.6	59.0
Δ_{PPO}	+3.3	+3.3	+2.5	+0.0	-0.1	+0.4	+0.5	+0.0	+2.2	+3.3	+1.1	+0.5	+0.0	-0.3	-0.2	+1.1
Mathstral-7B-v0.1	0.0	0.0	12.5	8.9	87.1	51.1	74.2	28.6	33.0	31.4	81.6	93.8	17.9	87.7	54.7	44.2
+ GRPO	0.0	0.0	47.5	17.8	92.9	55.9	81.3	35.7	44.0	49.2	88.1	97.6	25.6	93.0	81.5	54.0
+ GRPO w/ VERL.	6.7	0.0	45.0	20.0	93.3	55.5	81.5	50.0	40.7	46.6	89.5	97.2	29.3	90.7	83.5	55.3
Δ_{GRPO}	+6.7	+0.0	-2.5	+2.2	+0.4	-0.4	+0.2	+14.3	-3.3	-2.6	+1.4	-0.4	+3.7	-2.3	+2.0	+1.3
+ PPO	6.7	3.3	32.5	20.0	90.9	51.8	78.3	42.9	37.4	49.2	87.0	96.0	28.4	89.9	70.7	52.3
+ PPO w/ VERL.	10.0	0.0	27.5	22.2	93.0	53.8	78.2	42.9	51.6	48.3	87.4	96.7	26.1	89.6	84.1	54.1
Δ_{PPO}	+3.3	-3.3	-5.0	+2.2	+2.1	+2.0	-0.1	+0.0	+14.2	-0.9	+0.4	+0.7	-2.3	-0.3	+13.4	+1.8
Mistral-7B-v0.3	0.0	0.0	10.0	0.0	40.5	12.4	21.8	14.3	13.2	3.4	24.0	50.8	1.6	39.1	30.6	17.4
+ GRPO	0.0	0.0	2.5	4.4	58.2	11.1	42.3	0.0	15.4	5.1	52.4	79.2	3.0	47.6	37.7	23.9
+ GRPO w/ VERL.	0.0	0.0	7.5	2.2	59.1	15.0	43.0	0.0	6.6	4.2	40.3	69.5	2.8	57.5	53.0	24.0
Δ_{GRPO}	+0.0	+0.0	+5.0	-2.2	+0.9	+3.9	+0.7	+0.0	-8.8	-0.9	-12.1	-9.7	-0.2	+9.9	+15.3	+0.1
+ PPO	0.0	0.0	0.0	0.0	8.9	6.6	7.7	7.1	11.0	2.5	3.3	8.6	2.1	6.9	12.0	5.1
+ PPO w/ VERL.	0.0	0.0	2.5	0.0	44.7	10.6	35.7	7.1	16.5	5.1	28.8	70.7	2.4	57.5	35.1	21.1
Δ_{PPO}	+0.0	+0.0	+2.5	+0.0	+35.8	+4.0	+28.0	+0.0	+5.5	+2.6	+25.5	+62.1	+0.3	+50.6	+23.1	+16.0

G.5 Robustness Across Random Seeds

To assess robustness to random seed variation, we ran **five independent training seeds** for LLAMA-3.2-3B-INSTRUCT on all main benchmarks. We report mean \pm standard deviation over seeds in Tab. 8 (pass@1) and Tab. 9 (pass@k). Compared to GRPO, adding VERL improves the **average pass@1** from 39.24 ± 0.72 to 45.18 ± 0.38 (Tab. 8), and the **average pass@k** from 59.87 ± 1.37 to 62.78 ± 0.63 (Tab. 9). We further conduct two-sided **Welch’s t-tests** at significance level $\alpha = 0.05$ for each benchmark; statistically significant gains are marked with *. Overall, most improvements remain significant, suggesting the gains are not due to random seed noise.

Most benchmarks show statistically significant gains (marked with *). On the few tasks without *, the mean still improves, but variance across seeds is slightly larger. On these pass@k evaluations, GRPO+VERL again outperforms GRPO, with a statistically significant improvement in the average pass@k (see **Avg.** column).

G.6 Detailed Analysis of OOD Performance

A key motivation of this study is to assess whether our proposed VERL algorithm improves *out-of-domain (OOD)* generalization for reasoning-centric post-training, rather than merely overfitting to the in-domain training distribution. This is aligned with our overall claim that VERL enhances robustness and OOD performance beyond standard RL fine-tuning baselines.

To stress-test OOD generalization, we evaluate on a diverse suite of reasoning benchmarks spanning multiple domains: (i) **Astronomy**: AstroQA (English/Chinese) (Li et al., 2025a); (ii) **Science**: GPQA (Diamond/Extended/Main) (Rein et al., 2024) and SciBench (Wang et al., 2024b); (iii) **Law**: LEXam (Fan et al., 2025) and LexEval (Li et al., 2024); (iv) **Logical reasoning**: LogiQA (English/Chinese) (Liu et al., 2020) and ReClor (Yu et al., 2020); (v) **Medical QA**: MedMCQA (Pal et al., 2022), MedQA (Jin et al., 2021), and PubMedQA (Jin et al., 2019); (vi) **Broad knowledge**: MMLU (Hendrycks et al., 2021a) and MMLU-Pro (Wang et al., 2024c). All results are

Table 5: Performance comparison of instruction-tuned models under diverse decoding settings (Pass@ k). All results are reported in percentage (%).

Model	MATH500 (Pass@16)	AMC23 (Pass@128)	AMC24 (Pass@128)	AIME24 (Pass@256)	AIME25 (Pass@256)	Avg.
Llama-3.2-3B-Instruct	79.8	93.5	51.5	40.0	30.0	58.96
+ GRPO	80.2	95.4	60.6	40.0	30.0	61.24
+ GRPO w/ VERL.	80.6	95.7	59.0	50.0	36.7	64.40
Δ_{GRPO}	+0.4	+0.3	-1.6	+10.0	+6.7	+3.16
+ PPO	82.2	94.5	57.0	46.7	36.7	63.42
+ PPO w/ VERL.	82.4	94.7	57.8	46.7	40.0	64.32
Δ_{PPO}	+0.2	+0.2	+0.8	+0.0	+3.3	+0.90
Llama-3.1-8B-Instruct	79.8	94.6	57.4	46.7	36.7	63.04
+ GRPO	83.4	94.9	56.9	53.3	36.7	65.04
+ GRPO w/ VERL.	83.4	95.1	63.1	50.0	36.7	65.66
Δ_{GRPO}	+0.0	+0.2	+6.2	-3.3	+0.0	+0.62
+ PPO	79.2	92.4	59.0	46.7	36.7	62.80
+ PPO w/ VERL.	82.4	91.9	60.0	53.3	36.7	64.86
Δ_{PPO}	+3.2	-0.5	+1.0	+6.6	+0.0	+2.06
Qwen2.5-3B	86.0	96.7	69.0	56.7	40.0	69.68
+ GRPO	86.6	92.2	68.5	46.7	40.0	66.80
+ GRPO w/ VERL.	87.6	95.9	67.8	53.3	43.3	69.58
Δ_{GRPO}	+1.0	+3.7	-0.7	+6.6	+3.3	+2.78
+ PPO	87.8	96.5	67.9	43.3	43.3	67.76
+ PPO w/ VERL.	88.2	96.8	67.3	53.3	43.3	69.78
Δ_{PPO}	+0.4	+0.3	-0.6	+10.0	+0.0	+2.02
Qwen2.5-7B	90.6	98.4	73.7	60.0	60.0	76.54
+ GRPO	90.8	97.8	78.3	56.7	50.0	74.72
+ GRPO w/ VERL.	91.4	98.3	79.0	63.3	60.0	78.40
Δ_{GRPO}	+0.6	+0.5	+0.7	+6.6	+10.0	+3.68
+ PPO	91.2	98.6	74.3	53.3	56.7	74.82
+ PPO w/ VERL.	91.4	98.0	74.4	56.7	66.7	77.44
Δ_{PPO}	+0.2	-0.6	+0.1	+3.4	+10.0	+2.62
Mathstral-7B-v0.1	80.4	88.5	60.9	43.3	36.7	61.96
+ GRPO	84.8	87.3	69.2	36.7	40.0	63.60
+ GRPO w/ VERL.	87.0	97.0	76.9	50.0	50.0	72.18
Δ_{GRPO}	+2.2	+9.7	+7.7	+13.3	+10.0	+8.58
+ PPO	82.4	91.7	70.7	53.3	40.0	67.62
+ PPO w/ VERL.	84.8	93.8	69.9	53.3	46.7	69.70
Δ_{PPO}	+2.4	+2.1	-0.8	+0.0	+6.7	+2.08
Mistral-7B-v0.3	36.0	73.5	39.6	20.0	16.7	37.16
+ GRPO	33.0	63.2	36.0	10.0	10.0	30.44
+ GRPO w/ VERL.	34.4	64.5	38.0	16.7	13.3	33.38
Δ_{GRPO}	+1.4	+1.3	+2.0	+6.7	+3.3	+2.94
+ PPO	21.8	46.4	25.1	6.7	6.7	21.34
+ PPO w/ VERL.	19.2	46.5	30.1	3.3	13.3	22.48
Δ_{PPO}	-2.6	+0.1	+5.0	-3.4	+6.6	+1.14

Table 6: Comparison of GRPO + VERL using an intermediate layer (layer 14) versus the final layer, evaluated by pass@1 on multiple math benchmarks. Using the last layer yields the strongest average improvement.

Method	aime24	aime25	amc23	amc24	asdiv	carp_en	cmath	gaokao24_1	gaokao24_mix	gaokao_math_cloze	gsm8k	mawps	olympiadbench	svamp	tabmwp	Avg.
Llama-3.2-3B-Instruct	0.0	0.0	25.0	11.1	74.6	26.5	10.2	14.3	14.3	6.8	66.6	86.9	12.7	74.1	41.4	30.97
GRPO	3.3	0.0	27.5	8.9	88.8	45.0	28.3	21.4	20.9	23.7	80.7	96.0	16.7	87.7	71.7	41.37
GRPO w/ VERL (layer = 14)	10.0	0.0	27.5	11.1	88.6	43.6	30.7	21.4	16.5	22.0	81.9	95.5	18.1	87.0	71.4	41.69
GRPO w/ VERL (layer = last)	13.3	6.7	25.0	11.1	89.3	45.4	46.2	14.3	22.0	22.9	81.7	96.0	17.6	87.8	72.3	43.44

Table 7: Comparison of GRPO + VERL using an intermediate layer (layer 14) versus the final layer, evaluated by pass@k on several math benchmarks. Again, using the last layer yields the best average improvement.

Method	math500@16	amc23@128	amc24@128	aime24@256	aime25@256	Avg.
Llama-3.2-3B-Instruct	79.8	93.5	51.5	40.0	30.0	58.96
GRPO	80.2	95.4	60.6	40.0	30.0	61.24
GRPO w/ VERL (layer = 14)	81.0	92.9	57.0	40.0	36.7	61.52
GRPO w/ VERL (layer = last)	80.6	95.7	59.0	50.0	36.7	64.40

Table 8: Seed robustness on main benchmarks (pass@1). Results are mean \pm std over 5 seeds. * indicates Welch’s t-test significance at $\alpha = 0.05$ comparing GRPO+VERL against GRPO.

Method	aime24	aime25	amc23	amc24	asdiv	carp_cn	cmath	gaokao24_1	gaokao24_mix	gaokao_cloze	gsm8k	maaws	olympiad	svamp	tabmwp	Avg.
GRPO	10.00 \pm 0.00	0.00 \pm 0.00	17.50 \pm 0.00	11.54 \pm 0.98	88.62 \pm 0.54	44.82 \pm 0.04	15.34 \pm 0.36	5.72 \pm 7.83	19.58 \pm 2.95	22.20 \pm 1.82	80.88 \pm 0.60	95.60 \pm 0.00	17.80 \pm 0.85	87.08 \pm 0.22	71.92 \pm 0.41	39.24 \pm 0.72
GRPO+VERL	12.64 \pm 1.48*	3.98 \pm 1.52*	25.50 \pm 1.12*	16.00 \pm 2.46*	89.50 \pm 0.14*	45.20 \pm 0.14*	56.62 \pm 0.93*	25.72 \pm 3.94*	18.92 \pm 1.97	24.04 \pm 0.76	83.00 \pm 0.59*	96.32 \pm 0.16*	18.30 \pm 0.71	87.30 \pm 0.28	74.64 \pm 0.75*	45.18 \pm 0.38*

Table 9: Seed robustness on harder pass@k settings. Results are mean \pm std over 5 seeds. * indicates Welch’s t-test significance at $\alpha = 0.05$.

Method	math500@16	amc23@128	amc24@128	aime24@256	aime25@256	Avg.
GRPO	80.40 \pm 0.00	90.78 \pm 0.18	56.18 \pm 1.50	39.34 \pm 1.48	32.66 \pm 5.95	59.87 \pm 1.37
GRPO+VERL	80.52 \pm 0.18	94.22 \pm 1.01*	58.46 \pm 1.69	43.98 \pm 3.66*	36.70 \pm 0.00	62.78 \pm 0.63*

Table 10: Out-of-domain (OOD) reasoning performance (Pass@1, %) across diverse domains. Results are reported for Mathstral-7B-v0.1 and Qwen2.5-7B, comparing vanilla GRPO/PPO and their VERL-enhanced variants.

Model / Method	AstroQA EN	AstroQA ZH	GPQA D	GPQA Ext	GPQA Main	SciBench	LEXam	LexEval	LogiQA EN	LogiQA ZH	ReClor	MedMCQA	MedQA	PubMedQA	MMLU	MMLU-Pro	Avg.
Mathstral-7B-v0.1	50.7	43.8	29.3	26.0	25.4	15.8	19.6	38.7	45.8	44.6	43.6	42.9	39.7	22.4	53.9	26.4	35.25
+ GRPO	57.1	49.3	28.8	29.3	29.9	22.4	22.4	42.8	47.1	45.7	46.3	44.9	43.8	22.6	56.7	29.8	39.31
+ GRPO w/ VERL	58.0	51.6	31.8	31.1	34.6	27.9	22.6	47.6	50.4	44.8	49.7	47.1	46.1	24.4	58.8	30.7	41.68
Δ_{GRPO}	+0.9	+2.3	+3.0	+1.8	+4.7	+5.5	+0.2	+4.8	+3.3	-0.9	+3.4	+2.2	+2.3	+1.8	+2.1	+0.9	+2.37
+ PPO	52.7	47.5	28.8	29.9	27.9	22.1	17.3	44.5	47.5	45.0	46.5	44.2	45.2	23.1	56.5	29.1	37.93
+ PPO w/ VERL	54.5	49.1	31.3	26.6	26.1	23.0	21.4	42.2	47.5	51.6	46.2	44.8	45.1	22.3	56.9	29.3	39.34
Δ_{PPO}	+1.8	+1.6	+2.5	-3.3	-1.8	+0.9	+4.1	-2.3	+0.0	+6.6	-0.3	+0.6	-0.1	-0.8	+0.4	+0.2	+1.42
Qwen2.5-7B	61.2	58.0	27.3	26.6	24.6	15.8	21.1	50.4	50.8	49.7	56.0	50.6	62.0	43.6	61.3	36.1	43.44
+ GRPO	59.6	61.3	27.3	27.8	28.8	21.5	22.0	53.6	52.1	52.1	58.8	52.1	69.2	46.6	60.6	33.8	46.15
+ GRPO w/ VERL	59.4	62.4	30.8	28.9	26.8	22.8	23.0	54.6	52.8	52.4	61.6	52.0	70.3	47.8	60.9	33.4	47.09
Δ_{GRPO}	-0.2	+1.1	+3.5	+1.1	-2.0	+1.3	+1.0	+1.0	+0.7	+0.3	+2.8	-0.1	+1.1	+1.2	+0.3	-0.4	+0.94
+ PPO	60.0	62.1	27.3	25.3	28.3	22.1	22.8	55.2	52.1	52.1	58.6	53.6	70.4	45.8	60.9	34.1	46.20
+ PPO w/ VERL	60.9	61.3	34.8	30.0	27.0	23.0	23.4	57.1	53.0	52.3	61.8	52.9	71.0	45.4	61.0	34.4	47.45
Δ_{PPO}	+0.9	-0.8	+7.5	+4.7	-1.3	+0.9	+0.6	+1.9	+0.9	+0.2	+3.2	-0.7	+0.6	-0.4	+0.1	+0.3	+1.25

reported as Pass@1 (%).

Tab. 10 summarizes the OOD Pass@1 results for both Mathstral-7B and Qwen2.5-7B. Overall, VERL yields consistent average improvements when plugged into either GRPO or PPO: for Mathstral-7B, VERL improves the OOD average from 39.31 to 41.68 (+2.37) on GRPO and from 37.93 to 39.34 (+1.42) on PPO; for Qwen2.5-7B, it improves the OOD average from 46.15 to 47.09 (+0.94) on GRPO and from 46.20 to 47.45 (+1.25) on PPO. These gains are broadly distributed across heterogeneous domains, supporting that VERL enhances OOD reasoning robustness rather than trading off generalization for in-domain gains.

G.7 Comparison with Regularization Method in Action Space

This experiment aims to compare VERL with prior *action-space* regularization strategies used in RL post-training. Here, “action space” refers to the token-level policy over next-token actions during generation. A common practice is to regularize training using token-level proxy signals computed directly from this policy (e.g., confidence or entropy). Our goal is to test whether such action-

space proxy regularizers can match the performance gains of VERL, and to provide evidence that representation-level metrics yield a more effective regularization signal for reasoning.

Setup. We keep the post-training pipeline unchanged (same base model, RL algorithm, training data, optimizer, and decoding configuration), and only replace the auxiliary regularization signal. We consider three action-space baselines: (i) ACTION-LOGP, which regularizes using only the **average log-probability** of the generated response (a confidence-style signal); (ii) ACTION-ENTROPY, which regularizes using only the **response entropy** (a stochasticity-style signal); (iii) ACTION-LOGP+ENTROPY, which uses the **sum** of the above two signals. In contrast, VERL (VERL) employs hidden-state trajectory metrics defined in a semantically richer representation space. We report Pass@1 results in Tab. 11 and larger-budget decoding results (Pass@K) in Tab. 12.

Findings. Two observations stand out. First, action-space proxy regularization is often *unstable* under combination: while a single proxy can

Table 11: Comparison with action-space regularization baselines on selected reasoning datasets (Pass@1, %). Best value in each column is in **bold**; second best is underlined.

Model / Method	AIME24	AIME25	AMC23	AMC24	ASDiv	CARP_EN	CMATH	Gaokao24_I	Gaokao24_mix	Gaokao_math_cloze	GSM8K	MAWPS	OlympiadBench	SVAMP	TabMWP	Avg.
Llama-3.2-3B-Instruct	0.0	0.0	25.0	11.1	74.6	26.5	10.2	14.3	14.3	6.8	66.6	86.9	12.7	74.1	41.4	30.97
+ GRPO	3.3	0.0	<u>27.5</u>	8.9	88.8	45.0	28.3	21.4	20.9	23.7	80.7	96.0	16.7	87.7	71.7	41.37
+ GRPO w/ ACTION-LOGP regu	3.3	0.0	<u>27.5</u>	15.6	88.8	44.6	38.2	7.1	15.4	<u>32.0</u>	80.9	96.1	16.3	87.4	71.9	41.01
+ GRPO w/ ACTION-ENTROPY regu	<u>10.0</u>	0.0	35.0	8.9	89.2	<u>45.3</u>	37.5	7.1	16.5	28.0	81.9	96.1	17.9	87.0	71.6	<u>42.13</u>
+ GRPO w/ ACTION-LOGP+ENTROPY	3.3	3.3	25.0	<u>13.3</u>	88.5	44.7	37.7	0.0	18.7	<u>23.7</u>	82.0	96.1	15.6	87.0	70.5	40.63
+ GRPO w/ VERL (ours)	13.3	6.7	25.0	11.1	89.3	45.4	46.2	14.3	22.0	22.9	81.7	96.0	17.6	87.8	72.3	43.40

Table 12: Comparison with action-space regularization baselines under larger decoding budgets (Pass@K, %). Best value in each column is in **bold**; second best is underlined.

Model / Method	Math500 (Pass@16)	AMC23 (Pass@128)	AMC24 (Pass@128)	AIME24 (Pass@256)	AIME25 (Pass@256)	Avg.
Llama-3.2-3B-Instruct	79.8	93.5	51.5	40.0	<u>30.0</u>	58.96
+ GRPO	80.2	95.4	60.6	40.0	<u>30.0</u>	61.24
+ GRPO w/ ACTION-LOGP regu	<u>79.6</u>	93.8	56.9	40.0	<u>30.0</u>	60.06
+ GRPO w/ ACTION-ENTROPY regu	78.4	93.1	56.9	46.7	36.7	<u>62.36</u>
+ GRPO w/ ACTION-LOGP+ENTROPY	80.0	93.9	55.8	40.0	36.7	<u>61.28</u>
+ GRPO w/ VERL (ours)	80.6	95.7	<u>59.0</u>	50.0	36.7	64.40

sometimes help (e.g., ACTION-ENTROPY achieves the strongest or second-strongest performance on several columns), the naive combination ACTION-LOGP+ENTROPY does not reliably improve further and can even underperform the better single-signal choice (see the Avg. column in Tab. 11 and the Pass@K columns in Tab. 12). This indicates that token-level confidence (LogP) and stochasticity (entropy) may interfere when jointly used as a regularizer.

Second, across both evaluation regimes, VERL consistently provides the best overall results. On Pass@1, VERL achieves the highest Avg. and attains the best (or second-best) scores on most datasets in Tab. 11; on Pass@K, VERL similarly yields the strongest Avg. and the top performance on key high-budget reasoning benchmarks (e.g., AIME24 at Pass@256 and AMC23 at Pass@128 in Tab. 12). Overall, these results suggest that action-space token-level proxies are limited and do not reliably compose, whereas VERL offers a more robust and effective regularization signal for RL post-training.

G.8 Ablation on the Choice of Hidden Layer

We focus on the final hidden layer because our exploration/exploitation metrics are defined in the *semantic* space along the reasoning trajectory, and prior interpretability work (Jing et al., 2025; Sajjad et al., 2022; Valeriani et al., 2023; Matthews et al., 2024; Servedio et al., 2025; Zhang et al., 2025) suggests that the last layers are most aligned with semantic meaning and model predictions. Empirically, using the final layer gives consistently better performance than using an intermediate layer: for example, GRPO + VERL with the last layer im-

proves the average pass@1 from 41.69% (layer 14) to 43.44%, and the average pass@k from 61.52% (layer 14) to 64.40%. This subsection provides the detailed analysis supporting our design choice to base VERL on the final layer.

Intermediate layers in large language models can encode rich features. However, our notion of exploration and exploitation is explicitly defined in the semantic space of a reasoning trajectory. Existing interpretability studies (Jing et al., 2025; Sajjad et al., 2022; Valeriani et al., 2023; Matthews et al., 2024; Servedio et al., 2025; Zhang et al., 2025) indicate that hidden states in the last layers are most tightly aligned with token-level semantics and the model’s predictive distribution, while mid-layer representations tend to mix morphology, syntax, and other lower-level or task-specific signals. For this reason, we consider the final layer more suitable for semantic diversity metrics.

Intermediate-layer vs. final-layer VERL (pass@1). We use Llama-3.2-3B-Instruct as the base model, which has 28 transformer layers. We implemented VERL on an intermediate layer (layer 14) and on the last layer, keeping everything else fixed. The pass@1 results are reported in Tab. 6. Both VERL variants improve over GRPO, but the last-layer version clearly gives the strongest overall gains in average pass@1.

Intermediate-layer vs. final-layer VERL (pass@k). A similar pattern holds for pass@k, summarized in Tab. 7. Using the last layer yields the best average improvement, especially on the more challenging AIME-style benchmarks, further supporting the choice of the final hidden layer as the basis for our semantic exploration and

1956 exploitation metrics.

1957 H Case Study

1958 All case studies use Qwen2.5-7B as the base model.
1959 Outputs in the **gray boxes** are produced by vanilla
1960 GRPO trained for 120 steps, whereas outputs in
1961 the **purple boxes** are produced by VERL (GRPO
1962 + auxiliary shaping) trained for the same 120 steps.

1963 H.1 Case Study For Pass@1 Setting

1964 **Case Study I.** For the case in Fig. 11, vanilla
1965 GRPO incorrectly concludes that $-1 < -13$,
1966 while GRPO with our auxiliary shaping term pro-
1967 duces the correct comparison. This illustrates
1968 a weakness in vanilla GRPO’s exploitation: it
1969 does not reliably apply basic numerical reasoning,
1970 namely that among negative numbers, a larger ab-
1971 solute value implies a smaller number. Adding
1972 the auxiliary shaping term improves exploitation,
1973 making such comparisons more consistent.

1974 **Case Study II.** For the case in Fig. 12, the crucial
1975 distinction lies in how the models interpret the con-
1976 dition “*No house faces two different streets*”. The
1977 vanilla GRPO-trained model incorrectly infers that
1978 each street contains only 10 houses, overlooking
1979 the structural implication of the condition. In con-
1980 trast, the model trained with GRPO augmented by
1981 the auxiliary shaping term exploits this common-
1982 sense constraint more effectively: it decomposes
1983 the problem step by step, correctly deduces that
1984 each street has 20 houses, and thereby arrives at
1985 the correct total.

1986 **Case Study III.** As shown in Fig. 13, the task
1987 concerns solving a differential equation under the
1988 constraint of finding a purely exponential solution
1989 rather than a general one. The model trained with
1990 vanilla GRPO overlooks this constraint and defaults
1991 to the generic procedure for obtaining the full so-
1992 lution. In contrast, the model trained with GRPO
1993 augmented by the auxiliary shaping term respects
1994 the constraint and produces the correct purely ex-
1995 ponential answer. This case illustrates how the
1996 auxiliary shaping term enhances exploitation by en-
1997 forcing adherence to problem-specific constraints.

1998 H.2 Case Study For Pass@16 Setting

1999 **Case Study I.** As shown in Fig. 14, the question
2000 is a simple elementary-level math problem requir-
2001 ing the integer solution $x \geq 1, x \neq 2$. While the
2002 vanilla GRPO model consistently outputs 3, which
2003 is incorrect, the model trained with the auxiliary
2004 shaping term produces 1 in 5 out of 16 rollouts,
2005 the correct answer. This case highlights that the

2006 auxiliary shaping term enhances the model’s explo-
2007 ration ability, enabling it to consider a wider range
2008 of possibly valid solutions.

2009 **Case Study II.** As shown in Fig. 15, this high-
2010 school level geometry problem illustrates the con-
2011 trast between vanilla GRPO and GRPO with the
2012 auxiliary shaping term. The vanilla GRPO model
2013 fails to produce the correct answer in any roll-
2014 out. In contrast, the model trained with the auxil-
2015 iary shaping term succeeds in 5 out of 16 rollouts.
2016 Moreover, its reasoning traces display higher ex-
2017 ploratory diversity: it applies the law of sines 9
2018 times, the law of cosines 4 times, and other gen-
2019 eral geometric properties 3 times. By comparison,
2020 the vanilla GRPO model relies more narrowly on
2021 the law of sines (11 times), with only limited use
2022 of other strategies. This case highlights that the
2023 auxiliary shaping term enhances the model’s explo-
2024 ration capacity, enabling it to leverage a broader
2025 range of problem-solving strategies rather than
2026 over-committing to a single path.

Question

What is the smallest possible integer value of x if $12 \div (x + 1)$ is an integer?

Ground Truth

-13

Vanilla GRPO

To find the smallest possible integer value of x such that $\frac{12}{x+1}$ is an integer, we need to consider the factors of 12. The expression $\frac{12}{x+1}$ will be an integer if $x + 1$ is a divisor of 12. Let's list the divisors of 12:

1. The divisors of 12 are: $\pm 1, \pm 2, \pm 3, \pm 4, \pm 6, \pm 12$.
2. We need to find the smallest x , so we start with the smallest positive divisor, which is 1:

$$x + 1 = 1 \implies x = 1 - 1 \implies x = 0$$

3. To confirm, we substitute $x = 0$ back into the original expression:

$$\frac{12}{0 + 1} = \frac{12}{1} = 12$$

Since 12 is an integer, $x = 0$ is indeed a valid solution.

4. We should also consider the negative divisors to ensure we have the smallest possible integer value. The next smallest divisor is -1:

$$x + 1 = -1 \implies x = -1 - 1 \implies x = -2$$

5. Substituting $x = -2$ back into the original expression:

$$\frac{12}{-2 + 1} = \frac{12}{-1} = -12$$

Since -12 is an integer, $x = -2$ is also a valid solution, but it is not smaller than $x = 0$. Therefore, the smallest possible integer value of x is

0

GRPO with Auxiliary Shaping Term

To find the smallest possible integer value of x such that $\frac{12}{x+1}$ is an integer, we need to determine the divisors of 12. The divisors of 12 are $\pm 1, \pm 2, \pm 3, \pm 4, \pm 6, \pm 12$. For $\frac{12}{x+1}$ to be an integer, $x + 1$ must be one of these divisors. We want the smallest possible x , so we should consider the smallest positive divisor, which is 1.

If $x + 1 = 1$, then $x = 0$.

We should also consider the negative divisors to ensure we find the smallest possible integer value of x . The smallest negative divisor is -12. If

$x + 1 = -12$, then $x = -13$.

Comparing $x = 0$ and $x = -13$, the smallest possible integer value of x is -13.

Therefore, the final answer is

-13

Figure 11: Case study: Vanilla GRPO fails to exploit numerical commonsense when comparing negative values (-1 vs. -13), while our method gives the correct answer.

Question

In my neighborhood, there are six streets. There are 10 houses on each side of each street. No house faces two different streets. How many houses are in my neighborhood?

Ground Truth

120

Vanilla GRPO

To determine the total number of houses in the neighborhood, we can follow these steps:

1. **Identify the number of streets and houses per street:** - There are 6 streets. - Each street has 10 houses on each side.
2. **Calculate the number of houses on one street:** Since there are 10 houses on each side of a street, and each house faces only one street, the total number of houses on one street is 10 (since they are all on the same street).
3. **Calculate the total number of houses in the neighborhood:** There are 6 streets, and each street has 10 houses. Therefore, the total number of houses is $6 \times 10 = 60$.

Putting the final answer within , we get:

GRPO with Auxiliary Shaping Term

To find the total number of houses in the neighborhood, let's break down the problem step by step.

1. There are 6 streets in the neighborhood.
2. Each street has houses on both sides.
3. There are 10 houses on each side of every street.
4. No house faces two different streets, meaning each house is counted only once.

First, calculate the number of houses on one side of a single street:

$$10 \text{ houses/side}$$

Since there are houses on both sides of each street, multiply by 2:

$$10 \text{ houses/side} \times 2 \text{ sides/street} = 20 \text{ houses/street}$$

Now, since there are 6 streets, multiply the number of houses per street by the number of streets:

$$20 \text{ houses/street} \times 6 \text{ streets} = 120 \text{ houses}$$

So, the total number of houses in the neighborhood is

Figure 12: Case study: Vanilla GRPO overlooks the constraint and infers 10 houses per street, while our method exploits it correctly and infers 20.

Question

Find a purely exponential solution of $\frac{d^4x}{dt^4} - x = e^{-2t}$.

Ground Truth

$$\frac{e^{-2t}}{15}$$

Vanilla GRPO

To find a purely exponential solution of the differential equation ($\frac{d^4x}{dt^4} - x = e^{-2t}$), we will follow these steps:

1. **Find the complementary solution by solving the homogeneous equation:**

$$\frac{d^4x}{dt^4} - x = 0$$

The characteristic equation is:

$$r^4 - 1 = 0$$

Solving for r , we get:

$$r^4 = 1 \implies r = \pm 1, \pm i$$

Therefore, the complementary solution is:

$$x_c(t) = C_1 e^t + C_2 e^{-t} + C_3 \cos(t) + C_4 \sin(t)$$

2. **Find a particular solution using the method of undetermined coefficients:**

Since the right-hand side of the non-homogeneous equation is e^{-2t} , we assume a particular solution of the form:

$$x_p(t) = A e^{-2t}$$

Taking the fourth derivative of $x_p(t)$:

$$\frac{d^4}{dt^4}(A e^{-2t}) = A(-2)^4 e^{-2t} = 16A e^{-2t}$$

Substituting $x_p(t)$ and its fourth derivative into the original differential equation:

$$16A e^{-2t} - A e^{-2t} = e^{-2t}$$

Simplifying, we get:

$$(16A - A)e^{-2t} = e^{-2t} \implies 15A e^{-2t} = e^{-2t}$$

Therefore:

$$15A = 1 \implies A = \frac{1}{15}$$

So, the particular solution is:

$$x_p(t) = \frac{1}{15} e^{-2t}$$

3. **Combine the complementary and particular solutions:**

The general solution to the differential equation is:

$$x(t) = x_c(t) + x_p(t) = C_1 e^t + C_2 e^{-t} + C_3 \cos(t) + C_4 \sin(t) + \frac{1}{15} e^{-2t}$$

Thus, the final answer is:

$$C_1 e^t + C_2 e^{-t} + C_3 \cos(t) + C_4 \sin(t) + \frac{1}{15} e^{-2t}$$

Figure 13: Case study: Vanilla GRPO ignores the constraint of a purely exponential solution, while our method enforces it and yields the correct result.

GRPO with Auxiliary Shaping Term

To find a purely exponential solution to the differential equation

$$\frac{d^4 x}{dt^4} - x = e^{-2t}$$

, we will follow these steps:

1. **Find the complementary solution by solving the homogeneous equation:**

$$\frac{d^4 x}{dt^4} - x = 0$$

The characteristic equation is:

$$r^4 - 1 = 0$$

Solving for r , we get:

$$r^4 = 1 \implies r = \pm 1, \pm i$$

Therefore, the complementary solution is:

$$x_c(t) = C_1 e^t + C_2 e^{-t} + C_3 \cos(t) + C_4 \sin(t)$$

2. **Find a particular solution using the method of undetermined coefficients:** Since the right-hand side of the non-homogeneous equation is e^{-2t} , we assume a particular solution of the form:

$$x_p(t) = A e^{-2t}$$

Taking the fourth derivative of $x_p(t)$:

$$\frac{d^4 x_p}{dt^4} = 16A e^{-2t}$$

Substituting $x_p(t)$ and its fourth derivative into the original differential equation:

$$16A e^{-2t} - A e^{-2t} = e^{-2t}$$

Simplifying, we get:

$$15A e^{-2t} = e^{-2t}$$

Solving for A :

$$15A = 1 \implies A = \frac{1}{15}$$

Therefore, the particular solution is:

$$x_p(t) = \frac{1}{15} e^{-2t}$$

3. **Combine the complementary and particular solutions:** The general solution to the differential equation is:

$$x(t) = x_c(t) + x_p(t) = C_1 e^t + C_2 e^{-t} + C_3 \cos(t) + C_4 \sin(t) + \frac{1}{15} e^{-2t}$$

Thus, the purely exponential solution is:

$$\frac{1}{15} e^{-2t}$$

Figure 13: Case study: Vanilla GRPO ignores the constraint of a purely exponential solution, while our method enforces it and yields the correct result (continued).

Question

Given $f(x) = \frac{\sqrt{x-1}}{x-2}$, what is the smallest possible integer value for x such that $f(x)$ has a real number value?

Ground Truth

1

Figure 14: Case study demonstrating improved exploration: the auxiliary shaping term enables the model to explore and reach the correct solution.

(Details omitted as they are not crucial for this case study.)

1. **Numerator Condition:** The numerator is $\sqrt{x-1}$. For the square root to be defined and result in a real number, the expression inside the square root must be non-negative.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator is $x - 2$. For the function to be defined, the denominator must not be zero.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 3$. Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

To determine the smallest possible integer value for x such that $f(x) = \frac{\sqrt{x-1}}{x-2}$ has a real number value, we need to ensure that both the numerator and the denominator satisfy their respective conditions for real numbers.

1. **Numerator Condition:** The expression inside the square root must be non-negative.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator must not be zero to avoid division by zero.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these conditions, we need $x \geq 1$ and $x \neq 2$. The smallest integer that satisfies these conditions is $x = 3$. Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

To determine the smallest possible integer value for x such that $f(x) = \frac{\sqrt{x-1}}{x-2}$ has a real number value, we need to ensure that both the numerator and the denominator satisfy certain conditions.

1. **Numerator Condition:** The expression inside the square root must be non-negative because the square root of a negative number is not a real number.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator must not be zero because division by zero is undefined.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 3$.

(Details omitted as they are not crucial for this case study.)

1. **Square Root Condition:** The expression inside the square root, $x - 1$, must be greater than or equal to zero for the square root to be a real number.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator, $x - 2$, must not be zero to avoid division by zero.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 3$. Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

(Details omitted as they are not crucial for this case study.)

1. For the square root to be defined, the expression inside it must be greater than or equal to zero:

$$x - 1 \geq 0$$

Solving this inequality:

$$x \geq 1$$

2. For the denominator to be non-zero:

$$x - 2 \neq 0$$

Solving this equation:

$$x \neq 2$$

(Details omitted as they are not crucial for this case study.)

Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is:

$$\boxed{3}$$

(Details omitted as they are not crucial for this case study.)

1. For the square root to be defined, the expression inside it must be greater than or equal to zero:

$$x - 1 \geq 0$$

Solving this inequality:

$$x \geq 1$$

2. For the function to be defined, the denominator must not be zero:

$$x - 2 \neq 0$$

Solving this equation:

$$x \neq 2$$

(Details omitted as they are not crucial for this case study.)

Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

(Details omitted as they are not crucial for this case study.)

1. For the square root to be defined, we need:

$$x - 1 \geq 0$$

Solving this inequality:

$$x \geq 1$$

2. For the denominator to be non-zero, we need:

$$x - 2 \neq 0$$

Solving this equation:

$$x \neq 2$$

Combining these two conditions, x must be greater than or equal to 1 but not equal to 2. The smallest integer that satisfies this condition is $x = 3$. Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

(Details omitted as they are not crucial for this case study.)

1. **Square Root Condition:** The expression inside the square root, $x - 1$, must be greater than or equal to zero for the square root to be a real number.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator, $x - 2$, must not be zero to avoid division by zero.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 3$. Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

Figure 14: Case study demonstrating improved exploration: the auxiliary shaping term enables the model to explore and reach the correct solution (continued).

Vanilla GRPO#9-#16

(Details omitted as they are not crucial for this case study.)

1. **Numerator Condition:** The expression inside the square root, $x - 1$, must be non-negative because the square root of a negative number is not a real number. Therefore, we need:

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator, $x - 2$, must not be zero because division by zero is undefined. Therefore, we need:

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 3$. Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is: $\boxed{3}$

(Details omitted as they are not crucial for this case study.)

1. **Numerator Condition:** The expression inside the square root, $x - 1$, must be non-negative because the square root of a negative number is not a real number. Therefore, we need:

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator, $x - 2$, must not be zero because division by zero is undefined. Therefore, we need:

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 3$. Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is: $\boxed{3}$.

(Details omitted as they are not crucial for this case study.)

1. For the square root to be defined, we need:

$$x - 1 \geq 0$$

$$x \geq 1$$

2. For the denominator to be non-zero, we need:

$$x - 2 \neq 0$$

$$x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1$$

$$x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 3$. Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is: $\boxed{3}$.

(Details omitted as they are not crucial for this case study.)

1. **Numerator Condition:** The expression under the square root, $x - 1$, must be non-negative because the square root of a negative number is not a real number. Therefore, we need:

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator, $x - 2$, must not be zero because division by zero is undefined. Therefore, we need:

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 3$. Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is: $\boxed{3}$

(Details omitted as they are not crucial for this case study.)

1. **Numerator Condition:** The numerator is $\sqrt{x - 1}$. For the square root to be defined and result in a real number, the expression inside the square root must be non-negative.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator is $x - 2$. For the function to be defined, the denominator must not be zero.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we need x to be greater than or equal to 1, but not equal to 2. Therefore, the smallest possible integer value for x that satisfies both conditions is 3.

$\boxed{3}$

(Details omitted as they are not crucial for this case study.)

1. **Square Root Condition:** The expression inside the square root, $x - 1$, must be greater than or equal to zero for the square root to be a real number.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator, $x - 2$, must not be zero to avoid division by zero.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 3$. Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is: $\boxed{3}$.

To determine the smallest possible integer value for x such that $f(x) = \frac{\sqrt{x-1}}{x-2}$ has a real number value, we need to ensure that both the numerator and the denominator satisfy their respective conditions for real numbers.

1. **Numerator Condition:** The expression inside the square root must be non-negative.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator must not be zero.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these conditions, we need x to be greater than or equal to 1 but not equal to 2. Therefore, the smallest possible integer value for x that satisfies both conditions is 3.

$\boxed{3}$

(Details omitted as they are not crucial for this case study.)

1. **Numerator Condition:** The numerator is $\sqrt{x - 1}$. For the square root to be a real number, the expression inside the square root must be non-negative.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator is $x - 2$. For the function to be defined, the denominator must not be zero.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 3$. Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is: $\boxed{3}$.

Figure 14: Case study demonstrating improved exploration: the auxiliary shaping term enables the model to explore and reach the correct solution (continued).

GRPO with Auxiliary Shaping Term#1-#8

(Details omitted as they are not crucial for this case study.)

1. **Numerator Condition:** The numerator is $\sqrt{x-1}$. For the square root to be defined and real, the expression inside the square root must be non-negative.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator is $x - 2$. For the function to be defined, the denominator must not be zero.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these conditions, we need x to be greater than or equal to 1, but not equal to 2. The smallest integer that satisfies these conditions is $x = 3$.

Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

(Details omitted as they are not crucial for this case study.)

1. For the square root to be defined in the real numbers, the expression inside it must be greater than or equal to zero:

$$x - 1 \geq 0$$

Solving this inequality:

$$x \geq 1$$

2. For the function to be defined, the denominator must not be zero:

$$x - 2 \neq 0$$

Solving this equation:

$$x \neq 2$$

(Details omitted as they are not crucial for this case study.)

Thus, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

(Details omitted as they are not crucial for this case study.)

1. The expression inside the square root must be non-negative. 2. The denominator must not be zero.

(Details omitted as they are not crucial for this case study.)

Step 3: Combine the conditions From Step 1, we know $x \geq 1$. From Step 2, we know $x \neq 2$. Therefore, the smallest integer value for x that satisfies both conditions is the smallest integer greater than 1, which is 3.

Final Answer The smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

(Details omitted as they are not crucial for this case study.)

Combining the Conditions From Condition 1, we have $x \geq 1$. From Condition 2, we have $x \neq 2$.

Combining these, the smallest integer value for x that satisfies both conditions is $x = 3$.

Verification Let's verify that $x = 3$ works: - For $x = 3$, $\sqrt{x-1} = \sqrt{3-1} = \sqrt{2}$, which is a real number. - The denominator $x-2 = 3-2 = 1$, which is not zero.

Since both conditions are satisfied, the smallest possible integer value for x is indeed 3. Thus, the final answer is:

$$\boxed{3}$$

(Details omitted as they are not crucial for this case study.)

1. **Numerator Condition:** The expression under the square root, $x - 1$, must be non-negative because the square root of a negative number is not a real number. Therefore, we need:

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator, $x - 2$, must not be zero because division by zero is undefined. Therefore, we need:

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ but not $x = 2$ is $x = 1$.

Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{1}$.

(Details omitted as they are not crucial for this case study.)

Combining the Conditions From Condition 1, we know $x \geq 1$. From Condition 2, we know $x \neq 2$.

The smallest integer that satisfies $x \geq 1$ but not $x = 2$ is $x = 1$.

(Details omitted as they are not crucial for this case study.)

Since $f(1)$ is a real number, $x = 1$ is indeed a valid solution. Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{1}$.

(Details omitted as they are not crucial for this case study.)

1. For $\sqrt{x-1}$ to be defined, we need:

$$x - 1 \geq 0 \implies x \geq 1$$

2. For the denominator $x - 2$ to be non-zero, we need:

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 1$.

(Details omitted as they are not crucial for this case study.)

Since $f(1) = 0$ is a real number, our solution is confirmed. Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{1}$.

(Details omitted as they are not crucial for this case study.)

1. For $\sqrt{x-1}$ to be defined in the real numbers, we need:

$$x - 1 \geq 0 \implies x \geq 1$$

2. For the denominator $x - 2$ to not be zero, we need:

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, x must be greater than or equal to 1 but not equal to 2.

The smallest integer that satisfies this condition is $x = 1$.

Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{1}$.

Figure 14: Case study demonstrating improved exploration: the auxiliary shaping term enables the model to explore and reach the correct solution (continued).

GRPO with Auxiliary Shaping Term#9-#16

(Details omitted as they are not crucial for this case study.)

1. **Numerator Condition:** The expression inside the square root, $x - 1$, must be non-negative because the square root of a negative number is not a real number. Therefore, we need:

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator, $x - 2$, must not be zero because division by zero is undefined. Therefore, we need:

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 1$.

Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{1}$.

(Details omitted as they are not crucial for this case study.)

1. **Numerator Condition:** The expression under the square root must be non-negative.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator must not be zero.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 3$.

Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

(Details omitted as they are not crucial for this case study.)

1. **Numerator Condition:** The expression inside the square root must be non-negative.

$$x - 1 \geq 0 \implies x \geq 1$$

This ensures that $\sqrt{x-1}$ is a real number.

2. **Denominator Condition:** The denominator must not be zero.

$$x - 2 \neq 0 \implies x \neq 2$$

This ensures that the function is defined.

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 3$.

Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

(Details omitted as they are not crucial for this case study.)

1. For the square root to be defined, we need:

$$x - 1 \geq 0$$

Solving this inequality:

$$x \geq 1$$

2. For the denominator to be non-zero, we need:

$$x - 2 \neq 0$$

Solving this equation:

$$x \neq 2$$

Combining these two conditions, x must be greater than or equal to 1 but not equal to 2.

The smallest integer that satisfies this condition is $x = 3$.

Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

(Details omitted as they are not crucial for this case study.)

1. The expression inside the square root must be non-negative. 2. The denominator must not be zero.

Let's address these conditions one by one:

1. For the square root to be defined, the expression inside it must be greater than or equal to zero:

$$x - 1 \geq 0 \implies x \geq 1$$

2. For the denominator to be non-zero:

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these conditions, we need x to be greater than or equal to 1 but not equal to 2. The smallest integer that satisfies this is $x = 3$.

Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

(Details omitted as they are not crucial for this case study.)

1. **Numerator Condition:** The expression inside the square root must be non-negative because the square root of a negative number is not a real number.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator must not be zero because division by zero is undefined.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these two conditions, we get:

$$x \geq 1 \text{ and } x \neq 2$$

The smallest integer that satisfies $x \geq 1$ and $x \neq 2$ is $x = 3$.

Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

(Details omitted as they are not crucial for this case study.)

1. **Numerator Condition:** The expression under the square root, $x - 1$, must be non-negative because the square root of a negative number is not a real number.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator, $x - 2$, must not be zero because division by zero is undefined.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these conditions, we need $x \geq 1$ and $x \neq 2$. The smallest integer that satisfies these conditions is $x = 1$.

(Details omitted as they are not crucial for this case study.)

Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{1}$.

To determine the smallest possible integer value for x such that $f(x) = \frac{\sqrt{x-1}}{x-2}$ has a real number value, we need to ensure that both the numerator and the denominator satisfy their respective conditions for real numbers.

1. **Numerator Condition:** The expression inside the square root must be non-negative.

$$x - 1 \geq 0 \implies x \geq 1$$

2. **Denominator Condition:** The denominator must not be zero.

$$x - 2 \neq 0 \implies x \neq 2$$

Combining these conditions, we need x to be greater than or equal to 1 but not equal to 2. The smallest integer that satisfies this condition is $x = 3$.

Therefore, the smallest possible integer value for x such that $f(x)$ has a real number value is $\boxed{3}$.

Figure 14: Case study demonstrating improved exploration: the auxiliary shaping term enables the model to explore and reach the correct solution (continued).

Question

Abdul and Chiang are standing 48 feet apart in a field. Bharat is standing in the same field as far from Abdul as possible so that the angle formed by his lines of sight to Abdul and Chiang measures 60° . What is the square of the distance (in feet) between Abdul and Bharat?

Ground Truth

3072

Vanilla GRPO#1-#8

(Details omitted as they are not crucial for this case study.)

1. **Identify the triangle and its properties:**

(Details omitted as they are not crucial for this case study.)

2. **Use the Law of Sines:**
(Details omitted as they are not crucial for this case study.)

3. **Apply the Law of Sines:**
(Details omitted as they are not crucial for this case study.)

4. **Find the square of the distance:**

(Details omitted as they are not crucial for this case study.)

Therefore, the square of the distance between Abdul and Bharat is $\boxed{2304}$.

To solve this problem, we can use the Law of Sines in trigonometry. Let's denote the distance between Abdul and Bharat as x and the distance between Chiang and Bharat as y . We know that the angle between Abdul and Chiang as seen from Bharat is 60° , and the distance between Abdul and Chiang is 48 feet.

Given: - $\angle ABC = 60^\circ$ - $AB = 48$ feet - We need to find x^2

First, we apply the Law of Sines in triangle ABC :

(Details omitted as they are not crucial for this case study.)

Thus, the square of the distance between Abdul and Bharat is $\boxed{9216}$.

To solve this problem, we can use the Law of Sines in trigonometry. Let's denote the distance between Abdul and Bharat as AB , and the distance between Chiang and Bharat as BC . Since Bharat is as far from Abdul as possible to form a 60° angle, triangle ABC is a right triangle with the right angle at Chiang. This means $\angle ACB = 90^\circ$ and $\angle BAC = 60^\circ$.

(Details omitted as they are not crucial for this case study.)

Thus, the square of the distance between Abdul and Bharat is $\boxed{9216}$.

First, we need to find the distance AB using the Law of Sines. In triangle ABC :

(Details omitted as they are not crucial for this case study.)

Now, we need to find the square of the distance AB :

(Details omitted as they are not crucial for this case study.)

Thus, the square of the distance between Abdul and Bharat is $\boxed{1728}$.

To verify the accuracy of this result, we can use Python to perform the calculations.

(Details omitted as they are not crucial for this case study.)

To solve this problem, we can use the Law of Sines and some geometric properties. Let's break it down step by step.

1. **Identify the triangle and its properties:**

(Details omitted as they are not crucial for this case study.)

2. **Apply the Law of Sines:**
(Details omitted as they are not crucial for this case study.)

3. **Calculate the distances:**
(Details omitted as they are not crucial for this case study.)

4. **Find the square of the distance:**

(Details omitted as they are not crucial for this case study.)

So, the square of the distance between Abdul and Bharat is $\boxed{768}$.

(Details omitted as they are not crucial for this case study.)

1. Recognize that the triangle ABC is a right triangle with $\angle BAC = 90^\circ$ and $\angle ABC = 60^\circ$. Therefore, $\angle ACB = 30^\circ$.

(Details omitted as they are not crucial for this case study.)

3. Since $AC = 48$ feet, and in a 30-60-90 triangle, the side opposite the 60° angle is $\sqrt{3}$ times the side opposite the 30° angle, we can find AB :

$$AB = 48\sqrt{3}$$

(Details omitted as they are not crucial for this case study.)

Therefore, the square of the distance between Abdul and Bharat is $\boxed{6912}$ feet.

To solve this problem, we need to use some geometric principles, specifically involving the Law of Sines and properties of triangles.

1. **Identify the triangle and the given information:**

(Details omitted as they are not crucial for this case study.)

2. **Use the Law of Sines to find the distance BC (the distance from Bharat to Chiang):**

(Details omitted as they are not crucial for this case study.)

3. **Find the distance AC (the distance from Abdul to Bharat):**

(Details omitted as they are not crucial for this case study.)

Therefore, the square of the distance between Abdul and Bharat is $\boxed{11520}$.

(Details omitted as they are not crucial for this case study.)

1. **Identify the triangle and the given information:**

(Details omitted as they are not crucial for this case study.)

2. **Use the Law of Sines:**
(Details omitted as they are not crucial for this case study.)

3. **Maximize AB :**
(Details omitted as they are not crucial for this case study.)

4. **Apply the Pythagorean theorem:**
(Details omitted as they are not crucial for this case study.)

Thus, the square of the distance between Abdul and Bharat is $\boxed{9216}$.

Figure 15: Case study: Geometry problem. GRPO with auxiliary shaping term not only achieves higher accuracy but also explores more diverse reasoning strategies than vanilla GRPO.

Vanilla GRPO#9-#16

<p>To solve this problem, we can use the Law of Sines and some geometric properties. Let's break it down step by step.</p> <ol style="list-style-type: none"> Identify the triangle and the given information: (Details omitted as they are not crucial for this case study.) Determine the largest possible angle at A: (Details omitted as they are not crucial for this case study.) Use the Pythagorean theorem: (Details omitted as they are not crucial for this case study.) Find the square of the distance between Abdul and Bharat: (Details omitted as they are not crucial for this case study.) <p>So, the square of the distance between Abdul and Bharat is $\boxed{9216}$.</p>	<p>To solve this problem, we can use some geometric properties and trigonometry. Let's break it down step by step:</p> <ol style="list-style-type: none"> Identify the triangle formed: (Details omitted as they are not crucial for this case study.) Use the Law of Sines: (Details omitted as they are not crucial for this case study.) Maximize the distance: (Details omitted as they are not crucial for this case study.) Square the distance: (Details omitted as they are not crucial for this case study.) <p>Let's calculate it:</p> $48^2 = 2304$ <p>So, the square of the distance between Abdul and Bharat is $\boxed{2304}$.</p>	<p>To solve the problem, we need to determine the distance between Abdul and Bharat, given that the angle between their lines of sight to Chiang is 60° and that Bharat is positioned to be as far from Abdul as possible.</p> <p>(Details omitted as they are not crucial for this case study.)</p> <p>Using the Law of Cosines in triangle ABC:</p> <p>(Details omitted as they are not crucial for this case study.)</p> $BC^2 = d^2 + 2304 - 48d$ <p>(Details omitted as they are not crucial for this case study.)</p> <p>The square of the distance between Abdul and Bharat is:</p> $d^2 = 96^2 = 9216$ <p>Thus, the square of the distance between Abdul and Bharat is $\boxed{9216}$.</p>	<p>(Details omitted as they are not crucial for this case study.)</p> <ol style="list-style-type: none"> Identify the known values: (Details omitted as they are not crucial for this case study.) Set up the triangle: (Details omitted as they are not crucial for this case study.) Apply the Law of Sines: (Details omitted as they are not crucial for this case study.) Determine the angles: (Details omitted as they are not crucial for this case study.) Simplify the Law of Sines equation: (Details omitted as they are not crucial for this case study.) <p>Therefore, the square of the distance between Abdul and Bharat is $\boxed{2304}$.</p>
<p>To solve this problem, we need to use some geometric properties and trigonometry. Let's break it down step by step.</p> <ol style="list-style-type: none"> Understanding the Problem: (Details omitted as they are not crucial for this case study.) Geometric Insight: (Details omitted as they are not crucial for this case study.) Using Trigonometry: (Details omitted as they are not crucial for this case study.) Calculating the Distance: - Let the distance AB be x. Since $\angle ABC$ is 60° and we need to maximize AB, we use the Law of Cosines in triangle ABC: (Details omitted as they are not crucial for this case study.) <p>Therefore, the square of the distance between Abdul and Bharat is $\boxed{432}$.</p>	<p>(Details omitted as they are not crucial for this case study.)</p> <ol style="list-style-type: none"> Identify the given information: (Details omitted as they are not crucial for this case study.) Visualize the problem: (Details omitted as they are not crucial for this case study.) Use the Law of Sines: (Details omitted as they are not crucial for this case study.) Properties of an equilateral triangle: (Details omitted as they are not crucial for this case study.) Calculate the square of the distance: (Details omitted as they are not crucial for this case study.) <p>Thus, the square of the distance between Abdul and Bharat is $\boxed{2304}$.</p> <p>Providethefinalanswerwithin$\boxed{}$.</p>	<p>To solve this problem, we need to use some geometric properties and trigonometry. Let's break it down step by step.</p> <ol style="list-style-type: none"> Understand the Problem: (Details omitted as they are not crucial for this case study.) Visualize the Problem: (Details omitted as they are not crucial for this case study.) Use Geometric Properties: (Details omitted as they are not crucial for this case study.) Apply Trigonometry: (Details omitted as they are not crucial for this case study.) Calculate the Distance: - $AB = 2 * AC = 2 * 48 = 96$ feet. - The square of the distance AB is 96^2. Final Calculation: - $96^2 = 9216$. <p>So, the square of the distance between Abdul and Bharat is $\boxed{9216}$.</p>	<p>To solve this problem, we can use the Law of Sines in trigonometry. Let's break it down step by step:</p> <ol style="list-style-type: none"> Identify the given information: (Details omitted as they are not crucial for this case study.) Set up the triangle: (Details omitted as they are not crucial for this case study.) Use the Law of Sines: (Details omitted as they are not crucial for this case study.) Find the square of the distance: The distance between Abdul and Bharat is 24 feet. The square of this distance is: $24^2 = 576$ <p>Therefore, the square of the distance between Abdul and Bharat is $\boxed{576}$.</p>

Figure 15: Case study: Geometry problem. GRPO with auxiliary shaping term not only achieves higher accuracy but also explores more diverse reasoning strategies than vanilla GRPO (continued).

GRPO with Auxiliary Shaping Term#1-#8

(Details omitted as they are not crucial for this case study.)
 Given that $\angle ABC = 60^\circ$ and $AC = 48$ feet, we can use the Law of Sines to find AB :
 (Details omitted as they are not crucial for this case study.)
 The square of the distance AB is:

$$(32\sqrt{3})^2 = 32^2 \times 3 = 1024 \times 3 =$$

Thus, the square of the distance between Abdul and Bharat is 3072.

To solve this problem, we can use the Law of Sines in trigonometry.
 (Details omitted as they are not crucial for this case study.)
 Given that $AC = 48$ feet, we can use the Law of Sines to find AB :

$$\frac{AB}{\sin(60^\circ)} = \frac{AC}{\sin(30^\circ)}$$

(Details omitted as they are not crucial for this case study.)
 Solving for AB :

$$AB = 96 \times \frac{\sqrt{3}}{2} = 48\sqrt{3}$$

(Details omitted as they are not crucial for this case study.)
 Therefore, the square of the distance between Abdul and Bharat is 6912.

(Details omitted as they are not crucial for this case study.)

1. **Identify the given information:**

(Details omitted as they are not crucial for this case study.)
 2. **Visualize the problem:**
 (Details omitted as they are not crucial for this case study.)

3. **Determine the type of triangle:** - Since angle ABC and angle ACB are both 60° , triangle ABC is an equilateral triangle. In an equilateral triangle, all sides are equal.

4. **Calculate the side length of the equilateral triangle:** - Since $AC = 48$ feet and triangle ABC is equilateral, AB (the distance from Abdul to Bharat) is also 48 feet.

(Details omitted as they are not crucial for this case study.)
 So, the square of the distance between Abdul and Bharat is 2304.

To solve this problem, we can use the Law of Sines and some geometric properties. Let's denote the positions of Abdul, Chiang, and Bharat as points A , C , and B respectively. The distance between A and C is given as 48 feet. The angle $\angle ACB$ is 60° . We need to find the square of the distance AB , which is the farthest distance Bharat can be from Abdul.

(Details omitted as they are not crucial for this case study.)
 Therefore, the square of the distance between Abdul and Bharat is 6912.

To solve this problem, we need to use some geometric principles. Let's denote the positions of Abdul, Chiang, and Bharat as points A , C , and B respectively. We are given that the distance between A and C (AC) is 48 feet, and the angle ABC is 60° . We need to find the distance between A and B (AB), and then square it.
 (Details omitted as they are not crucial for this case study.)
 Step 3: Recognize that if triangle ABC is equilateral, then all sides are equal. Therefore, $AB = BC = AC = 48$ feet.
 Step 4: Calculate the square of the distance AB . Since $AB = 48$ feet, we square this value to get AB^2 .

$$AB^2 = 48^2 = 2304$$

So, the square of the distance between Abdul and Bharat is 2304.

To find the square of the distance between Abdul and Bharat, we can use the Law of Sines in trigonometry. Let's denote the distance between Abdul and Bharat as AB , and the distance between Chiang and Bharat as BC . Since Bharat is as far from Abdul as possible, triangle ABC is a right triangle with $\angle ACB = 90^\circ$.
 Given: - $\angle ABC = 60^\circ$ - $AC = 48$ feet
 We want to find AB^2 .
 First, we apply the Law of Sines in triangle ABC :
 (Details omitted as they are not crucial for this case study.)
 Thus, the square of the distance between Abdul and Bharat is 3072.

(Details omitted as they are not crucial for this case study.)

In a right triangle, if one of the angles is 60° , the triangle is a 30-60-90 triangle. In a 30-60-90 triangle, the sides are in the ratio $1 : \sqrt{3} : 2$. Here, AC is the hypotenuse, so we can set up the relationship as follows:

$$AC = 2 \times AB$$

$$48 = 2 \times AB$$

$$AB = \frac{48}{2} = 24 \text{ feet}$$

We need to find the square of the distance AB :

$$AB^2 = 24^2 = 576$$

Thus, the square of the distance between Abdul and Bharat is 576.

(Details omitted as they are not crucial for this case study.)

Using the properties of a 30-60-90 triangle, we know that the side opposite the 60° angle is $\sqrt{3}$ times the side opposite the 30° angle. Let's denote the side opposite the 30° angle (which is half the distance between Abdul and Chiang) as x . Therefore, $AB = x\sqrt{3}$ and $BC = 2x$.

(Details omitted as they are not crucial for this case study.)

$$AB^2 = 2304(21 - 12\sqrt{3})$$

Since we are asked to find the square of the distance, and we have already calculated it, we can box the final answer:

$$\boxed{2304}$$

Figure 15: Case study: Geometry problem. GRPO with auxiliary shaping term not only achieves higher accuracy but also explores more diverse reasoning strategies than vanilla GRPO (continued).

GRPO with Auxiliary Shaping Term#9-#16

<p>To solve this problem, we can use the Law of Sines in trigonometry. (Details omitted as they are not crucial for this case study.) Using the Law of Sines:</p> $\frac{AB}{\sin(120^\circ)} = \frac{AC}{\sin(60^\circ)}$ <p>(Details omitted as they are not crucial for this case study.) Therefore, the square of the distance between Abdul and Bharat is $\boxed{3072}$.</p>	<p>To solve this problem, we need to use some geometric properties and trigonometric relationships. Let's break it down step by step.</p> <ol style="list-style-type: none"> Identify the given information: (Details omitted as they are not crucial for this case study.) Visualize the problem: (Details omitted as they are not crucial for this case study.) Use the Law of Cosines: (Details omitted as they are not crucial for this case study.) Apply the Law of Cosines: (Details omitted as they are not crucial for this case study.) Find the square of the distance: - The square of the distance between Abdul and Bharat is $d^2 = 2304$. Therefore, the square of the distance between Abdul and Bharat is $\boxed{2304}$. 	<p>To solve this problem, we can use the Law of Sines and some geometric properties. (Details omitted as they are not crucial for this case study.) Given that the angle ABC is 60° and AB is maximized, the triangle ABC is a right triangle with angle ACB being 90° (since the angle at Bharat is 60° and the sum of angles in a triangle is 180°, the remaining angle at A must be 30°, making AB the hypotenuse of a 30-60-90 triangle). (Details omitted as they are not crucial for this case study.) Therefore, the square of the distance between Abdul and Bharat is $\boxed{3072}$.</p>	<p>To solve this problem, we need to use some geometric principles. Let's denote the positions of Abdul, Chiang, and Bharat as points A, C, and B respectively. The distance between A and C is 48 feet, and we need to find the distance AB when the angle $\angle ABC$ is 60° and B is as far from A as possible.</p> <ol style="list-style-type: none"> Identify the configuration: (Details omitted as they are not crucial for this case study.) Use the Law of Cosines: (Details omitted as they are not crucial for this case study.) Substitute into the Law of Cosines: (Details omitted as they are not crucial for this case study.) So, the square of the distance between Abdul and Bharat is $\boxed{3840}$ square feet.
<p>To solve this problem, we need to use some geometric principles, specifically involving the Law of Sines and properties of triangles. Let's break it down step by step:</p> <ol style="list-style-type: none"> Identify the given information: (Details omitted as they are not crucial for this case study.) Visualize the problem: (Details omitted as they are not crucial for this case study.) Use the Law of Sines: (Details omitted as they are not crucial for this case study.) Determine the angles: (Details omitted as they are not crucial for this case study.) Calculate the distance AB: (Details omitted as they are not crucial for this case study.) Therefore, the square of the distance between Abdul and Bharat is $\boxed{3072}$. 	<p>To find the square of the distance between Abdul and Bharat, let's denote the distance between Abdul and Bharat as d, and the distance between Chiang and Bharat as d'. Since the angle formed by the lines of sight to Abdul and Chiang is 60°, we can use the Law of Cosines in the triangle formed by Abdul, Chiang, and Bharat. The Law of Cosines states that for any triangle with sides a, b, and c, and the angle θ opposite side c:</p> $c^2 = a^2 + b^2 - 2ab \cos(\theta)$ <p>(Details omitted as they are not crucial for this case study.) Therefore, the square of the distance between Abdul and Bharat is $\boxed{2304}$.</p>	<p>To solve this problem, we'll use some geometric principles, particularly the Law of Sines and properties of triangles. Let's break down the problem step-by-step:</p> <ol style="list-style-type: none"> Identify the given information: (Details omitted as they are not crucial for this case study.) Visualize the problem: (Details omitted as they are not crucial for this case study.) Use the Law of Sines: (Details omitted as they are not crucial for this case study.) Apply the Law of Sines: (Details omitted as they are not crucial for this case study.) Use trigonometry in the right triangle: (Details omitted as they are not crucial for this case study.) Therefore, the square of the distance between Abdul and Bharat is $\boxed{1728}$. 	<p>To solve this problem, we need to use some geometric properties and trigonometry. Let's denote the positions of Abdul, Chiang, and Bharat as points A, C, and B respectively in the field. We know that AC = 48 feet and the angle $\angle ACB = 60^\circ$. We want to find the distance AB, which we'll denote as d, and then find d^2.</p> <ol style="list-style-type: none"> Identify the type of triangle: (Details omitted as they are not crucial for this case study.) Apply the Law of Cosines: (Details omitted as they are not crucial for this case study.) Maximize the distance d: (Details omitted as they are not crucial for this case study.) Therefore, the square of the distance between Abdul and Bharat is $\boxed{1728}$.

Figure 15: Case study: Geometry problem. GRPO with auxiliary shaping term not only achieves higher accuracy but also explores more diverse reasoning strategies than vanilla GRPO (continued).