
RAGEN-2: Reasoning Collapse in Agentic RL

Anonymous Authors¹

Abstract

RL training of multi-turn LLM agents is unstable, and reasoning quality drives task performance. Entropy, the standard reasoning-stability monitor, only measures within-input diversity and misses whether reasoning depends on the input. We identify **template collapse**: stable entropy alongside input-agnostic boilerplate, invisible to entropy and existing metrics. We diagnose it via a **mutual-information (MI) proxy** that scores cross-input distinguishability online; across tasks, MI correlates with final performance far more strongly than entropy. We then explain collapse via a **signal-to-noise ratio (SNR)** mechanism: low within-input reward variance weakens task gradients, letting input-agnostic regularization dominate and erase cross-input differences. We mitigate this with **SNR-Aware Filtering**, prioritizing high-variance prompts each iteration. Across planning, math reasoning, web navigation, and code execution, the method consistently improves input dependence and task performance.

1. Introduction

Training multi-turn LLM agents with reinforcement learning (RL) is inherently challenging (Qi et al., 2025; Zhang et al., 2026; Yu et al., 2025). Researchers therefore monitor reward for **outcome stability** and entropy for **reasoning process stability** (Schulman et al., 2017b; Ouyang et al., 2022; Xu et al., 2025), treating both as stability indicators of RL training.

However, entropy can be an ambiguous signal to understand reasoning quality. When entropy decreases, it may simply reflect the model becoming more specialized and confident on the task, which is a natural outcome of RL optimization (Yu et al., 2025; Xu et al., 2025). When entropy remains high, reasoning can still drift toward fixed

templates that appear diverse within any single input but are effectively the same across inputs (Figure 1). We call this **template collapse**, a failure mode invisible to both metrics. This risk is especially acute in multi-turn settings: sparse rewards cannot distinguish input-driven reasoning from templated reasoning that merely happens to succeed (Wang & Ammanabrolu, 2025; Wang et al., 2025c), and reasoning chains are hard to get directly supervised (Shao et al., 2024; Cui et al., 2025). As a result, template collapse can persist unnoticed during training, making agents unreliable and silently hurting their reasoning abilities.

To understand and mitigate template collapse, this paper addresses two questions. **(Q1) How to diagnose?** (§2) Entropy-based metrics (Wei et al., 2025; Yao et al., 2025; Yun et al., 2025) track within-input variability but miss input dependence across inputs, so they fail to detect template collapse. We propose a mutual information (MI) proxy (Cover & Thomas, 2006) that scores each reasoning chain against all batch inputs to measure input dependence, without external models. **(Q2) Why does it happen?** (§3) We explain through a signal-to-noise ratio (SNR) lens. Task gradients draw signal from reward differences across within-input trajectories. Sampling noise and input-agnostic regularization (KL divergence and entropy regularization (Schulman et al., 2017b; Xu et al., 2025)) dilute this signal. Low SNR lets noise dominate, erasing cross-input reasoning differences.

To address template collapse, based on the SNR view, we introduce **SNR-Aware Filtering**, which uses reward variance as a lightweight SNR proxy to select high-signal prompts each iteration, without additional supervision. Throughout training, the MI proxy monitors input dependence; across experiments, MI correlates with task performance significantly more strongly than entropy, validating it as a diagnostic for template collapse.

Together, they constitute a diagnostic framework for a systematic failure mode in multi-turn agent RL, validated across planning (Schrader, 2018), mathematical reasoning (Yu et al., 2023; Katz et al., 2025), web navigation, code execution, and tool use, under multiple RL algorithms, model scales, and modalities. SNR-Aware Filtering consistently improves input dependence and task performance, providing direct experimental support for the SNR mechanism.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

Our contributions are summarized as follows:

1. **Identifying template collapse.** We find that template collapse occurs when reasoning appears diverse within inputs but becomes input-agnostic across inputs. We propose a mutual information proxy to detect it without external models.
2. **Explaining template collapse via SNR.** We show that low reward variance weakens task gradients while input-agnostic regularization remains constant, erasing input dependence. We provide gradient decomposition evidence across reward-variance buckets.
3. **SNR-Aware Filtering.** We propose filtering prompts by reward variance before each update. We demonstrate that this improves input dependence and performance across tasks, algorithms, scales, and modalities.

2. Template Collapse in Multi-turn Agent RL

2.1. Setup and Preliminaries

We study closed-loop multi-turn agent reinforcement learning (Wang et al., 2025c), where a policy π_θ is trained by repeatedly rolling out trajectories under the current policy and environment and updating on the collected experience. At each time step t , the agent observes o_t , generates a response consisting of reasoning tokens z_t and an executable action a_t , and receives reward r_t , forming a trajectory $\tau = \{(o_t, z_t, a_t, r_t)\}_{t=1}^T$.

We use X to denote the full context available to the model immediately before generating reasoning at turn t : this comprises the system prompt, all prior observations $o_{1:t}$, actions $a_{1:t-1}$, and reasoning tokens $z_{1:t-1}$. We use Z to denote the reasoning token sequence the model generates for that turn, excluding action tokens and boundary markers (e.g., `</think>`).

The standard PPO/GRPO objective contains regularization terms (KL divergence, entropy bonus) that act uniformly across all inputs regardless of their content:

$$\mathcal{L}(\theta) = \mathbb{E}_{x,\tau} [A(\tau, x)] - \lambda_{\text{KL}} D_{\text{KL}}(\pi_\theta \| \pi_{\text{ref}}) + \lambda_H H(\pi_\theta),$$

where $A(\tau, x)$ is the advantage.

2.2. Rethinking Reasoning Collapse from an Information-Theoretic Lens

Why entropy is insufficient to measure reasoning quality? Researchers proxy process stability with entropy and outcome stability with reward, treating both as evidence of healthy training. Stable entropy, however, does not guarantee stable reasoning. Reasoning diversity (marginal entropy) $H(Z)$ decomposes via Cover & Thomas (2006):

$$H(Z) = I(X; Z) + H(Z | X), \quad (1)$$

where $I(X; Z)$ is input dependence (**mutual information** between input X and reasoning Z), and $H(Z | X)$ is within-

input diversity (**conditional entropy** of reasoning given input). Entropy metrics proxy $H(Z | X)$, but neither captures a decline in $I(X; Z)$: the policy can sustain high $H(Z | X)$ while $I(X; Z)$ drops to zero, producing diverse but input-agnostic boilerplate. We call this **template collapse**.

Reasoning regimes with a mutual information view. Figure 1 illustrates four reasoning states along these two axes: (i) *Diverse Reasoning* (high $H(Z | X)$, high $I(X; Z)$): the desired regime where reasoning is both varied within each input and systematically grounded across different inputs; (ii) *Template Collapse* (high $H(Z | X)$, low $I(X; Z)$): superficially diverse but input-agnostic—the systematic blind spot of existing stability metrics; (iii) *Compressed Reasoning* (low $H(Z | X)$, high $I(X; Z)$): input-faithful but overly deterministic; and (iv) *Low-Entropy Collapse* (low $H(Z | X)$, low $I(X; Z)$): fully degenerate with deterministic and input-agnostic outputs. Among these, Template Collapse is uniquely problematic because entropy-based metrics can remain high while input dependence collapses. Empirically, $I(X; Z)$ correlates significantly more strongly with task performance than entropy does (Figure 9).

2.3. Mutual Information Proxy Family

How do we estimate mutual information? True mutual information $I(X; Z)$ has no closed form for high-dimensional token sequences, so we propose an empirical proxy $\hat{I}(X; Z)$ based on retrieval. The intuition: mutual information $I(X; Z)$ measures how much knowing the reasoning Z tells us about which input X produced it. When $I(X; Z)$ is high, different inputs yield distinguishable reasoning patterns—the model adapts its reasoning to the specific problem. When $I(X; Z)$ is low, reasoning becomes input-agnostic: observing Z gives little clue about which X it came from. This is the signature of template collapse. If reasoning truly collapses into templates, it should be easy to detect: a reasoning trace Z generated from input X_i will be equally likely under any other input X_j .

Method: In-Batch Cross-Scoring. Given P prompts and G reasoning samples per prompt from training rollouts, we compute teacher-forced log-likelihoods for every $(Z_{i,k}, X_j)$ pair, forming the scoring matrix $\mathbf{L}_{i,k,j} = \log p_\theta(Z_{i,k} | X_j)$. We extract two length-normalized quantities:

$$\begin{aligned} \text{matched}_{i,k} &= \frac{\mathbf{L}_{i,k,i}}{|Z_{i,k}|}, \\ \text{marginal}_{i,k} &= \frac{1}{|Z_{i,k}|} \log \frac{1}{P} \sum_j \exp(\mathbf{L}_{i,k,j}). \end{aligned} \quad (2)$$

where $\text{matched}_{i,k}$ is the per-token log-likelihood of reasoning $Z_{i,k}$ under its true source input X_i , and $\text{marginal}_{i,k}$ approximates the marginal log-likelihood $\log p_\theta(Z_{i,k})$ via a uniform mixture over all prompts in the batch.

Two Primary Proxies. We use two complementary proxies

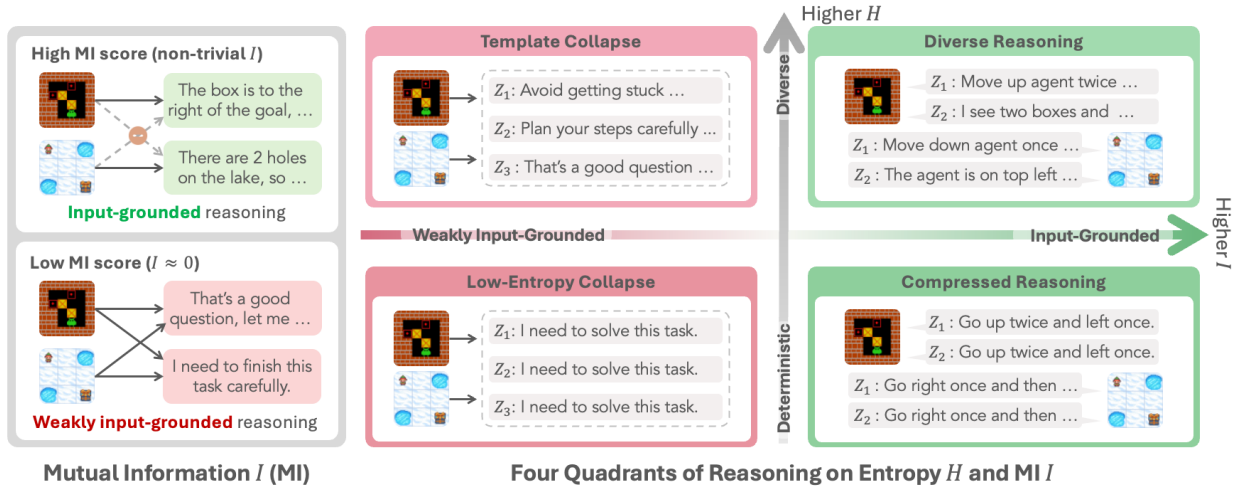


Figure 1. Left: input-driven reasoning adapts to the current state; templated reasoning produces nearly identical responses across different inputs. Right: four reasoning regimes characterized along two axes: conditional entropy $H(Z | X)$ (within-input diversity) and mutual information $I(X; Z)$ (input dependence). Details in Section 2.

derived from Eq. 2:

(1) *Retrieval-Acc* (discrete, interpretable): We define

$$\text{Acc} = \frac{1}{PG} \sum_{i=1}^P \sum_{k=1}^G \mathbb{I} \left[i = \arg \max_j \mathbf{L}_{i,k,j} \right].$$

Under collapse, Acc approaches chance level $1/P$ (1.56% at $P=64$), providing an absolute reference.

(2) *MI-ZScore-EMA* (continuous, robust): We estimate input dependence as

$$\hat{I}(X; Z) = \frac{1}{PG} \sum_{i=1}^P \sum_{k=1}^G \left(\text{matched}_{i,k} - \text{marginal}_{i,k} \right),$$

which increases when reasoning is much more compatible with its source input than with the batch mixture. In template-collapse regimes, $\text{matched}_{i,k} \approx \text{marginal}_{i,k}$ for many samples and thus $\hat{I}(X; Z)$ approaches 0. We apply z-score normalization and exponential moving average (EMA) to stabilize training monitoring, yielding MI-ZScore-EMA.

Proxy Variants and Validation. Appendix B lists additional proxy variants, varying along three dimensions: (1) turn scope (first-turn only vs. trajectory-uniform sampling); (2) aggregation (discrete retrieval vs. continuous MI estimate); (3) length normalization (per-token vs. per-sequence). For comparison, conditional entropy $H(Z | X) = -\frac{1}{PG} \sum_{i,k} \text{matched}_{i,k}$ and marginal entropy $H(Z) = -\frac{1}{PG} \sum_{i,k} \text{marginal}_{i,k}$ are logged in parallel, satisfying $H(Z) = \hat{I}(X; Z) + H(Z | X)$. We set $\epsilon = 10^{-3}$ and $\alpha = 0.9$ for z-score normalization and EMA, respectively.

Empirically, Retrieval-Acc and MI-ZScore-EMA achieve positive Spearman correlation with final task performance

(+0.39 for Trajectory MI-ZScore), substantially above entropy metrics, which show negative correlations (-0.11 to -0.14), confirming entropy is misleading in direction (Figure 9). All proxies reuse $(X_i, Z_{i,k})$ pairs from the training rollout and require no additional model or inference pass; implementation details are in Appendix E.

3. The Mechanism of Template Collapse: A Signal-to-Noise Ratio (SNR) View

We have defined template collapse (low $I(X; Z)$, high $H(Z | X)$) and introduced an MI proxy to diagnose it. This section explains why RL training produces this failure mode and how to mitigate it. Our core finding: when policy gradient updates are dominated by input-agnostic noise rather than task-discriminative signal—low signal-to-noise ratio (SNR)—reasoning drifts toward templates that appear diverse within each input but ignore cross-input differences.

3.1. Observing Signal-Noise Imbalance in RL Gradients

We begin with an empirical observation that motivates the mechanistic analysis. Sorting training prompts by their within-input reward variance $\widehat{\text{Var}}(R | X)$ and grouping them into equal-sized buckets, we measure the gradient norms contributed by task objectives / regularization terms (Figure 3). Three patterns are consistent across algorithms:

- Task gradient scales with reward variance:** $\|g_{\text{task}}\|$ increases monotonically with bucket RV. High-variance prompts yield strong task-discriminative gradients; low-variance prompts produce weak gradients even when non-zero.
- Regularization gradient is flat:** $\|g_{\text{reg}}\|$ (from KL and

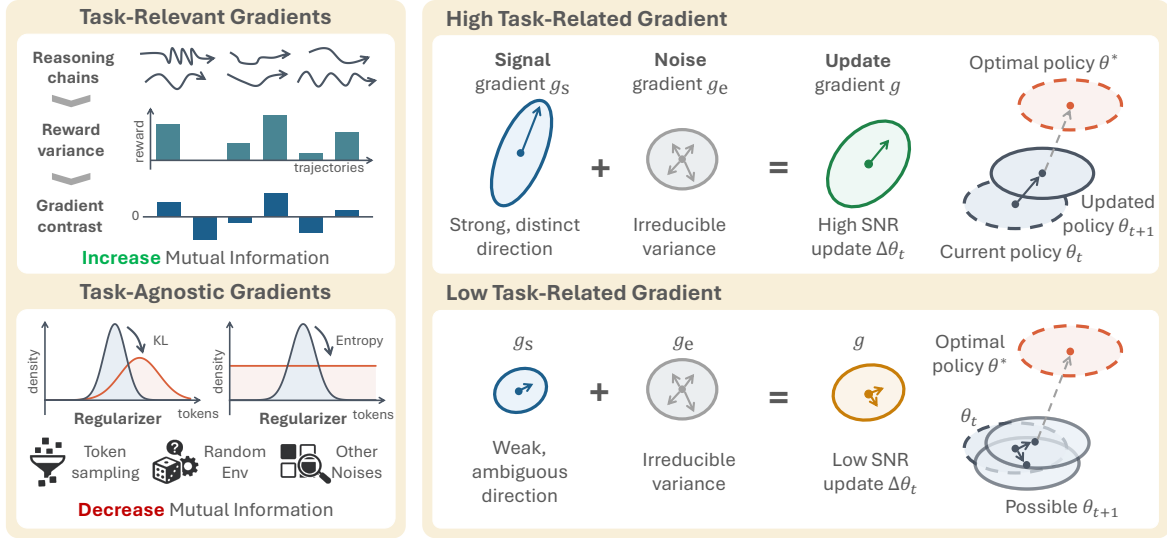


Figure 2. Schematic Signal-to-Noise Ratio (SNR) view of RL updates. Left: total gradient decomposes into task gradient (sharpenes with higher within-input reward variance) and regularization gradient. Right: high reward variance yields strong task gradient and better convergence (high SNR); low reward variance makes regularization gradient dominate, producing erratic updates and input-agnostic reasoning (low SNR).

entropy terms) remains constant across all buckets, applying uniform contraction to every reasoning chain regardless of its source prompt or reward signal.

- Low-RV prompts produce gradient updates dominated by regularization:** In the lowest-variance buckets, task gradients nearly vanish while regularization gradients persist, meaning updates are driven almost entirely by input-agnostic noise.

This gradient imbalance suggests that low reward variance weakens the task-discriminative component of updates, allowing input-agnostic regularization to dominate. When many prompts fall into this regime, the model learns to produce reasoning that satisfies regularization constraints (diverse, fluent) but ignores input-specific requirements—exactly the signature of template collapse.

3.2. Formalizing the SNR Mechanism via Gradient Decomposition

The empirical pattern above can be formalized through a *signal-to-noise decomposition of policy gradients*. Low within-input reward variance collapses advantages toward zero, weakening the task gradient. Simultaneously, input-agnostic regularization terms apply uniform contraction to every reasoning chain regardless of its source prompt. When the task gradient is weak, regularization dominates every update and pushes reasoning toward input-agnostic patterns, lowering $I(X; Z)$. This is the gradient-level mechanism behind template collapse (Figure 2; regularizer-dominance analysis in Appendix L).

For input x with G sampled trajectories, the advantage esti-

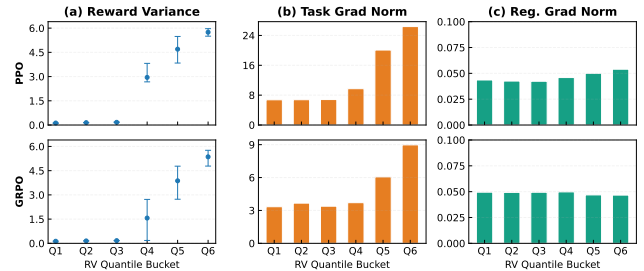


Figure 3. Prompts sorted into six reward-variance buckets Q1–Q6: (a) task gradient norm rises monotonically with bucket RV; (b) low-RV task gradients persist but carry almost no useful signal; (c) regularizer gradient norm (KL + entropy) is flat. Supports the SNR mechanism under both algorithms.

mate is $A_g = R_g - \bar{R}(x)$ and the task gradient is

$$g_{\text{task}}(x) = \frac{1}{G} \sum_g A_g \nabla_{\theta} \log \pi_{\theta}(\tau_g | x).$$

The Cauchy-Schwarz inequality gives (Appendix D):

$$|g_{\text{task}}(x)| \leq \sqrt{\widehat{\text{Var}}(R | X = x)} \cdot C.$$

Low reward variance therefore weakens g_{task} while leaving g_{reg} unchanged, driving $I(X; Z) \rightarrow 0$. Critically, $H(Z | X)$ need not decline: entropy regularization can sustain within-input diversity while input dependence collapses.

We formalize this through a three-noise decomposition of the total gradient:

$$g_{\text{total}} = g_{\text{signal}} + g_{\text{task-noise}} + g_{\text{reg}}.$$

Signal and task noise both vary across prompts, but only the former carries task-discriminative information. Regularization noise acts uniformly at the chain level: every reasoning chain receives the same KL/entropy contraction regardless

Table 1. Three-noise decomposition of the policy update gradient.

Component	Source	Mitigation
g_{signal}	Reward differences across same-prompt trajectories	SNR-Aware Filtering
$g_{\text{task-noise}}$	Sampling and environment stochasticity	Filter high-noise prompts
g_{reg}	Uniform per-chain contraction (KL, entropy)	Tune $\lambda_{\text{KL}}, \lambda_{\text{ent}}$

of its source prompt, making it inherently input-agnostic and the direct suppressive force on cross-input differences (Table 1).

In practice, $g_{\text{task}} = g_{\text{signal}} + g_{\text{task-noise}}$ merges the two prompt-level components, and the SNR is

$$\text{SNR}(x) = \frac{\|g_{\text{signal}}(x)\|}{\|g_{\text{task-noise}}(x)\| + \|g_{\text{reg}}\|}.$$

Low SNR shifts updates toward input-agnostic directions, lowering $I(X; Z)$ even when $H(Z | X)$ remains high (Appendix L).

3.3. SNR-Aware Filtering: Prioritizing High-Signal Updates

The gradient analysis above identifies the *mechanism behind template collapse*: low reward variance weakens task signal, allowing regularization noise to dominate and push reasoning toward input-agnostic patterns. This suggests a direct mitigation strategy: prioritize prompts with higher within-input reward variance, where advantage estimates carry stronger task-discriminative information and regularization is less likely to dominate the update.

We propose **SNR-Aware Filtering**: at each training iteration, estimate $\widehat{\text{Var}}(R | X)$ for each prompt and retain only the top fraction by variance before computing parameter updates (workflow in Figure 4). This concentrates gradient budget on high-SNR prompts and filters out low-variance updates that would be dominated by input-agnostic regularization.

Reward variance as SNR proxy. At each iteration, we estimate $\text{Var}(R | X)$ at the prompt level by sampling G trajectories for the same prompt X and computing the sample variance of episode returns:

$$\widehat{\text{Var}}(R | X) = \frac{1}{G-1} \sum_{g=1}^G (R_g(X) - \bar{R}(X))^2,$$

$$\bar{R}(X) = \frac{1}{G} \sum_{g=1}^G R_g(X).$$

Higher $\widehat{\text{Var}}(R | X)$ indicates trajectories can be meaningfully distinguished by reward, strengthening advantage estimates and increasing the likelihood that gradients align with task-relevant directions (Appendix I).

Table 2. Summary of the features of the environments used.

Task	Stochastic	Multi-turn	State	Reward
Sokoban	✗	✓	Grid	Dense
FrozenLake	✓	✓	Grid	Binary
MetaMathQA	✗	✓	Text	Dense
Countdown	✗	✗	Text	Binary
SearchQA	✗	✓	Text	Dense
WebShop	✗	✓	Text	Dense
DeepCoder	✗	✗	Text	Dense

Top- p filtering by reward variance. We keep the top fraction of prompts by variance score with keep rate $\rho \in (0, 1]$, analogous to nucleus sampling (Holtzman et al., 2020) but ranking by per-prompt reward variance. Let $V_i = \widehat{\text{Var}}(R | X = x_i)$ and let σ be a permutation that sorts prompts by descending V , so $V_{\sigma(1)} \geq \dots \geq V_{\sigma(P)}$. With threshold $\tau = \rho \sum_i V_i$, the kept set is $S = \{\sigma(1), \dots, \sigma(k^*)\}$ where $k^* = \min\{k : \sum_{j=1}^k V_{\sigma(j)} \geq \tau\}$, and the filtered objective is $\mathcal{L}_\rho(\theta) = \frac{1}{k^*} \sum_{i \in S} \sum_{j \in \mathcal{B}_i} L_\theta(\xi_j)$. This concentrates updates on high-signal prompts while adapting the kept count to the variance distribution. Other filtering strategies (top- k , min- p) and implementation details are in Appendix H.

4. Experiments

We first establish that template collapse occurs reliably across training configurations (Section 4.2), then evaluate SNR-Aware Filtering as an intervention across tasks, algorithms, model scales, and modalities (Section 4.3; Table 3).

4.1. Experimental Testbed

We adopt the RAGEN (Wang et al., 2025c) testbed and evaluate LLM agents on four controllable tasks that stress complementary decision-making regimes: irreversible planning (Sokoban), sparse-reward long-horizon navigation under stochastic transitions (FrozenLake), and symbolic math reasoning (MetaMathQA, Countdown). To further evaluate multi-turn reasoning and decision-making capabilities, we also include SearchQA (Tan et al., 2025), WebShop (Yao et al., 2022), and DeepCoder (Mattern et al., 2025; Li et al., 2023; Jain et al., 2024) (see Appendix C.1 for detailed descriptions).

Training and evaluation setup. We train Qwen2.5-3B (Qwen Team, 2024) with the veRL/HybridFlow stack (Sheng et al., 2024), following RAGEN (Wang et al., 2025c) defaults unless otherwise stated. We compare PPO (Schulman et al., 2017b), DAPO (Yu et al., 2025), GRPO (Shao et al., 2024), and Dr. GRPO (Liu et al., 2025) for up to 400 rollout-update iterations. Each iteration collects $K = P \times G = 128$ trajectories per environment, with prompt batch size $P = 8$ and group size $G = 16$ trajectories per prompt. When applying SNR-Aware Filtering with

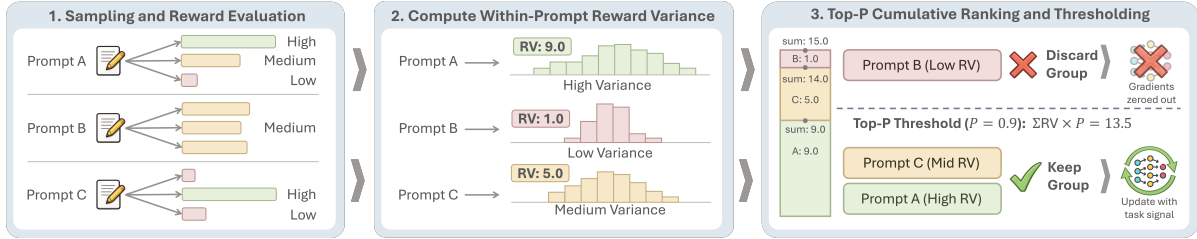


Figure 4. SNR-Aware Filtering workflow. At each training iteration: (1) rollout generation collects trajectories; (2) within-prompt reward variance is computed as SNR proxy; (3) prompts are ranked by RV and top- p fraction retained; policy update is performed only on the high-signal subset. The loop prevents updates on noisy rollouts and requires no additional models or rollouts beyond standard RL.

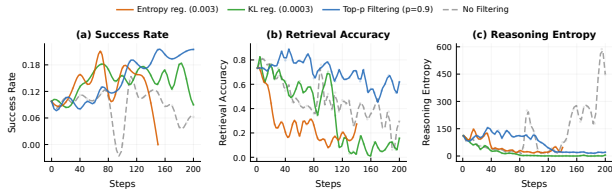


Figure 5. Training dynamics under different intervention strategies. (a) Task success rate, (b) MI proxy (retrieval accuracy), and (c) reasoning entropy. Without filtering, MI degrades early while entropy spikes, signaling template collapse. Filtering effectively mitigates the decline in retrieval accuracy, with top- p SNR-Aware filtering best preserving both task performance and reasoning diversity.

keep rate ρ , we reduce the effective minibatch size accordingly and scale the per-step loss by ρ , so the optimization step size remains comparable.

4.2. Template Collapse as a Consistent Failure Mode

Across all training configurations, RL-trained agents reliably develop reasoning that is fluent but input-agnostic: $I(X; Z)$ declines while $H(Z | X)$ remains high, and this drift is invisible to entropy-based monitoring.

Observing template collapse through MI dynamics. We track three key metrics during training: task success rate, our MI proxy $\hat{I}(X; Z)$ (Retrieval-Acc), and conditional entropy $H(Z | X)$ (Figure 5). We present dynamics for all MI proxies in Appendix E.1. The trajectory reveals a critical pattern: mutual information declines significantly before task performance degrades, while conditional entropy remains elevated throughout. This divergence is the hallmark of template collapse. Reasoning appears diverse within each input (high $H(Z | X)$) but becomes increasingly input-agnostic across inputs (low $I(X; Z)$).

The early decline of $\hat{I}(X; Z)$ demonstrates that our MI proxy serves as an early warning signal, detecting reasoning degradation that entropy-based metrics miss entirely. This finding motivates using MI as a primary diagnostic alongside task performance, rather than relying solely on entropy for process monitoring.

Behavioral manifestation of template collapse. Reasoning length declines monotonically across eight environments spanning spatial (Yin et al., 2025), logic puzzle (Chen et al.,

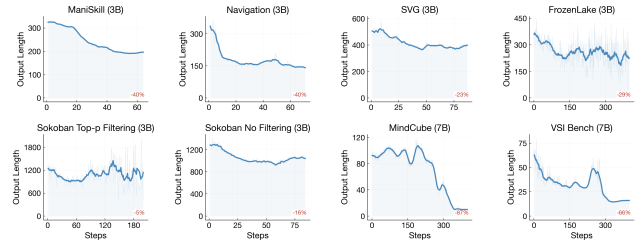


Figure 6. Reasoning length decline across eight environments, showing systematic compression as a behavioral signature of template collapse.

2025), visual (Wang et al., 2025b), and math (Liu et al., 2026) agents (Figure 6), a behavioral signature that complements MI-based diagnostics.

4.3. SNR-Aware Filtering Consistently Improves Performance

Comparing filtering strategies. Top- p (nucleus-style) filtering consistently outperforms Top- k (fixed-count) and no-filter baselines across four environments (Figure 7). We use Top- p as the default in all subsequent experiments.

Table 3 summarizes our experimental matrix over four tasks, multiple RL algorithms, model scales/types, and input modalities. Across this grid, SNR-Aware Filtering yields two consistent effects. First, it improves peak task success rate in most settings (reported as the $+\Delta$ next to each peak), demonstrating that prioritizing high-signal updates strengthens learning efficiency. Second, the gains span multiple experimental axes, including (i) the RL optimizer (PPO / DAPO / GRPO / Dr. GRPO), (ii) the model family and scale (Qwen2.5 from 0.5B to 7B; Llama3.2-3B), and (iii) the input modality (text- and image-conditioned Qwen2.5-VL). Here, DAPO and Dr. GRPO are recent strong baselines that directly target stable training and mitigate collapse-like failure modes. In Table 3, DAPO “no-filter” results correspond to the original algorithms without our filtering applied. DAPO itself also includes a filtering/acceptance step; it can be interpreted as a special case of our framework where the selection is fixed (equivalently, a top- P filter with $P \rightarrow 1.0$), while our SNR-Aware Filtering provides an explicit, tunable SNR knob via the keep rate ρ . This breadth suggests SNR-Aware Filtering serves as a general-purpose SNR control

Table 3. SNR-Aware Filtering results (%) across algorithms, model scales, types, and modalities. Each cell reports baseline peak with filter delta in parentheses; Qwen2.5-VL-3B includes text (T) and image (V) inputs. Filtering improves average score across all variants.

Experiment Variants	Sokoban	FrozenLake	MetaMathQA	Countdown	Average
Baseline					
PPO (Schulman et al., 2017b), Qwen2.5-3B (Qwen Team, 2024)	12.9 (+16.0)	67.0 (+10.9)	92.6 (+0.6)	97.9 (+0.0)	67.6 (+6.9)
Algorithm					
DAPO (Yu et al., 2025)	16.2 (+5.1)	66.8 (+2.1)	90.8 (+2.8)	95.7 (+1.6)	67.4 (+2.9)
GRPO (Shao et al., 2024)	12.1 (+9.0)	70.9 (-3.0)	91.2 (+1.2)	95.7 (+2.2)	67.5 (+3.7)
Dr. GRPO (Liu et al., 2025)	12.1 (-0.4)	23.2 (+0.6)	91.2 (+1.4)	96.5 (+1.4)	55.8 (+0.8)
Model Scale (PPO)					
Qwen2.5-0.5B (Qwen Team, 2024)	3.3 (+22.9)	19.5 (+0.0)	10.0 (-0.2)	23.0 (-0.7)	14.0 (+5.5)
Qwen2.5-1.5B (Qwen Team, 2024)	17.0 (+6.2)	36.5 (+1.6)	80.3 (+7.0)	56.6 (+1.6)	47.6 (+4.1)
Qwen2.5-7B (Qwen Team, 2024)	42.4 (+4.9)	85.0 (-0.6)	84.0 (+11.7)	97.7 (+0.3)	77.3 (+4.1)
Model Type					
Qwen2.5-3B-Instruct (Qwen Team, 2024)	22.5 (+14.2)	83.6 (+2.3)	91.2 (+0.4)	96.3 (-0.6)	73.4 (+4.1)
Llama3.2-3B (Meta Llama, 2024)	24.4 (+18.8)	84.6 (-0.2)	86.1 (+3.7)	99.2 (-1.2)	73.6 (+5.3)
Modality (Input Type)					
Qwen2.5-VL-3B (T) (Bai et al., 2025)	53.0 (+6.0)	16.0 (+53.5)	-	-	34.5 (+29.8)
Qwen2.5-VL-3B (V) (Bai et al., 2025)	65.0 (+12.0)	19.5 (+59.5)	-	-	42.3 (+35.8)

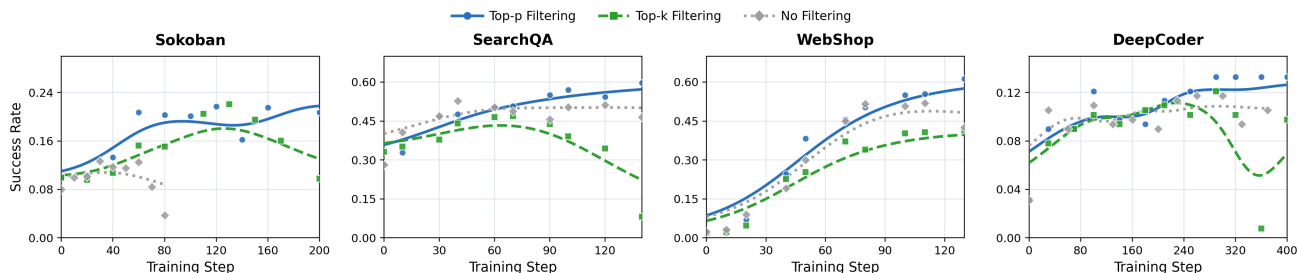


Figure 7. Top- p filtering consistently outperforms Top- k and no-filter baselines across four environments.

knob and works alongside standard stabilization terms (e.g., KL and entropy regularization).

Compute overhead. SNR-Aware Filtering needs $G \geq 2$ trajectories per prompt to estimate per-prompt RV; the variance computation itself adds $< 0.1\%$ of iteration time, and filtering reduces per-step gradient time by 26–41% at $\rho = 0.9$. Full $P \times G$ sweep on Sokoban in Appendix R.1.

5. Analysis

5.1. MI Diagnoses Collapse Better Than Entropy Across All Interventions

We demonstrate that MI separates high- and low-performance runs across all three intervention families better, and entropy could conflate them. At the same training budget, stronger SNR-Aware Filtering moves runs toward higher MI and better performance; KL and entropy tuning shift entropy without moving MI. We sweep three families of interventions (entropy regularization strength, KL constraint strength, and SNR-Aware Filtering keep rate) and compare their trajectories in both diagnostic spaces at fixed training steps (Figure 8). Entropy- and KL-based stabilizers

induce larger changes in $H(Z | X)$ than in $\hat{I}(X; Z)$, and rarely move the model into the high- $\hat{I}(X; Z)$ regime with clearly improved performance. In contrast, SNR-Aware Filtering traces a monotone improvement in both $\hat{I}(X; Z)$ and task success; pushing entropy too high leads to instability and performance collapse, while KL constraint mainly anchors the policy near its reference distribution without boosting input dependence.

We compute Spearman correlation between task success rate and each candidate diagnostic across runs with varying entropy regularization strength, KL constraint strength, and Top- p filtering kept mass (Figure 9). MI-family metrics achieve positive correlations, with Trajectory MI-ZScore reaching +0.39. In contrast, Reasoning Entropy and Conditional Entropy metrics show near-zero or negative correlations (between -0.11 and -0.14). This confirms that MI predicts performance twice as reliably as entropy does, and entropy actually points in the wrong direction. These results validate MI as a superior training monitor compared to entropy-based diagnostics for multi-turn agent RL.

Stress-testing the SNR claim. The SNR account makes concrete causal predictions about how reward variance, envi-

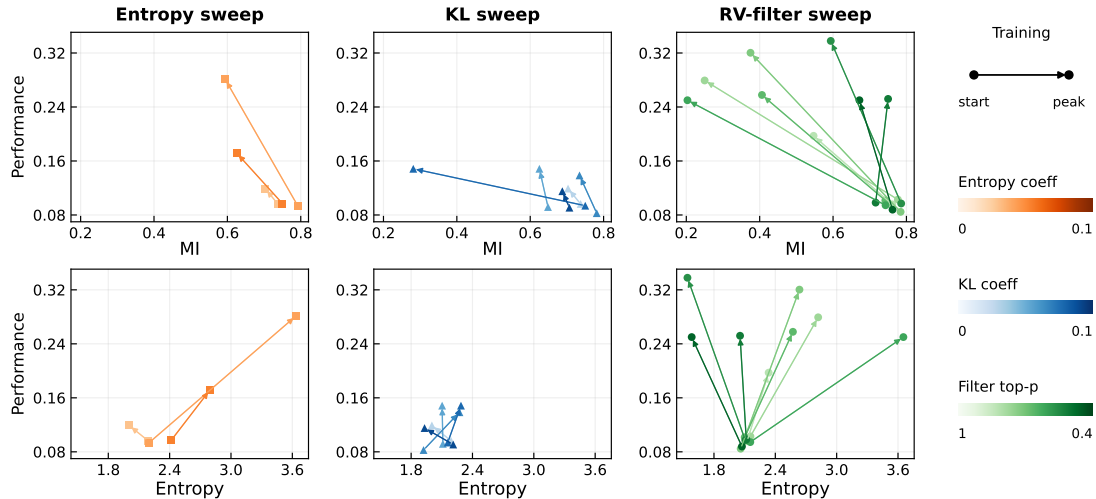


Figure 8. Training dynamics under three interventions. For each setting we connect two checkpoints (steps 10/400) into a trajectory (arrows point to later steps); color intensity is weaker to stronger intervention.

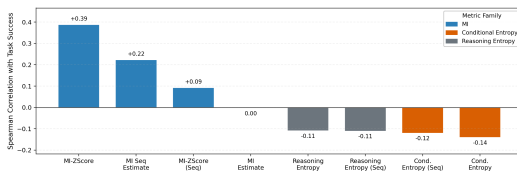


Figure 9. Spearman correlations showing MI-family metrics more positively predict performance than entropy metrics.

ronmental noise, and prompt selection should affect $\hat{I}(X; Z)$ and task performance. Appendix S verifies these via quartile ablation, controlled noise injection, prompt- vs. trajectory-level filtering, a Std/Mean(RV) applicability diagnostic, the adaptive kept-ratio dynamics, and a comparison against KL/entropy tuning.

6. Related Work

Reasoning collapse and policy degeneracy. LLM-agent RL reports reasoning collapse (templated rationales) and policy-level degeneracy (behavior concentrating on easy-to-reproduce patterns) (DeepSeek AI, 2025; Wei et al., 2025; Yao et al., 2025; Yun et al., 2025; Feng et al., 2025; Wang & Ammanabrolu, 2025), echoing model collapse in self-training even when average metrics look stable (Gerstgrasser et al., 2024; Shumailov et al., 2024).

Evaluating reasoning diversity and input dependence. Most diversity metrics — lexical (Li et al., 2016; Zhu et al., 2018), embedding-based (Pillutla et al., 2021; Tevet & Berant, 2021), uncertainty (Montahaei et al., 2019; Semeniuta et al., 2019) — capture within-input variability without testing whether differences are input-driven (Tevet & Berant, 2021; Yun et al., 2025); recent input-dependence probes include behavioral tests (Gardner et al., 2020; Ribeiro et al., 2020; Zhu et al., 2024) and retrieval-style matching (Morris et al., 2023; Gao et al., 2024; Zhang et al., 2024; Li & Klabjan, 2025). Reasoning faithfulness (Lanham et al., 2023;

Turpin et al., 2023; Siegel et al., 2024; Zaman & Srivastava, 2025) asks a different question (whether rationales reflect true decision bases).

Stabilizing closed-loop Agent RL. Prior work spans KL/entropy control, clipping, reward shaping, and curricula (Schulman et al., 2017a;b; Haarnoja et al., 2019; Ouyang et al., 2022; Rafailov et al., 2024; Feng et al., 2025; Wang & Ammanabrolu, 2025; Xu et al., 2025; Yao et al., 2025), plus stepwise rewards and self-correction for multi-step agents (Cobbe et al., 2021; Uesato et al., 2022; Madaan et al., 2023; Shinn et al., 2023; Yao et al., 2023; Wei et al., 2025); none target the within-input reward-variance axis we identify.

7. Conclusions and Limitations

We find closed-loop multi-turn agent RL can fail silently: reasoning drifts toward fluent but input-agnostic boilerplate while conditional entropy remains stable. We define this as **template collapse**. Built on this, the paper makes three contributions. First, we introduce a mutual information (MI) proxy between inputs and reasoning, which interprets template collapse and tracks task performance better than conditional entropy. To explain why collapse occurs, we propose SNR mechanism in RL and show that low within-input reward variance suppresses task gradients and lets regularization forces dominate, pushing policy outputs toward input-agnostic templates. To address this, we introduce SNR-Aware Filtering to prioritize prompts with reward variance before each parameter update, improving performance on average across tasks, model scales, and modalities and can integrate easily with existing training pipelines.

Limitations. The SNR decomposition assumes signal/noise separability; the method needs RV as a reliable signal proxy and degrades in sparse or highly stochastic reward regimes; all experiments are single-agent; aggressive filtering can narrow exploration. Detailed discussion in Appendix T.

References

- Bai, S., Chen, K., Liu, X., Wang, J., Ge, W., Song, S., Dang, K., Wang, P., Wang, S., Tang, J., Zhong, H., Zhu, Y., Yang, M., Li, Z., Wan, J., Wang, P., Ding, W., Fu, Z., Xu, Y., Ye, J., Zhang, X., Xie, T., Cheng, Z., Zhang, H., Yang, Z., Xu, H., and Lin, J. Qwen2.5-vl technical report, 2025. URL <https://arxiv.org/abs/2502.13923>.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. Openai gym, 2016. URL <https://arxiv.org/abs/1606.01540>.
- Chen, S., Zhu, T., Wang, Z., Zhang, J., Wang, K., Gao, S., Xiao, T., Teh, Y. W., He, J., and Li, M. Internalizing world models via self-play finetuning for agentic rl, 2025. URL <https://arxiv.org/abs/2510.15047>.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>.
- Cover, T. M. and Thomas, J. A. *Elements of Information Theory*. Wiley-Interscience, 2 edition, 2006.
- Cui, G., Yuan, L., Wang, Z., Wang, H., Zhang, Y., Chen, J., Li, W., He, B., Fan, Y., Yu, T., Xu, Q., Chen, W., Yuan, J., Chen, H., Zhang, K., Lv, X., Wang, S., Yao, Y., Han, X., Peng, H., Cheng, Y., Liu, Z., Sun, M., Zhou, B., and Ding, N. Process reinforcement through implicit rewards, 2025. URL <https://arxiv.org/abs/2502.01456>.
- DeepSeek AI. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081): 633–638, September 2025. ISSN 1476-4687. doi: 10.1038/s41586-025-09422-z. URL <http://dx.doi.org/10.1038/s41586-025-09422-z>.
- Dou, Z.-Y., Yang, C.-F., Wu, X., Chang, K.-W., and Peng, N. Re-rest: Reflection-reinforced self-training for language agents, 2025. URL <https://arxiv.org/abs/2406.01495>.
- Feng, L., Xue, Z., Liu, T., and An, B. Group-in-group policy optimization for llm agent training, 2025. URL <https://arxiv.org/abs/2505.10978>.
- Gao, L., Peng, R., Zhang, Y., and Zhao, J. Dory: Deliberative prompt recovery for llm, 2024. URL <https://arxiv.org/abs/2405.20657>.
- Gardner, M., Artzi, Y., Basmova, V., Berant, J., Bogin, B., Chen, S., Dasigi, P., Dua, D., Elazar, Y., Gottumukkala, A., Gupta, N., Hajishirzi, H., Ilharco, G., Khashabi, D., Lin, K., Liu, J., Liu, N. F., Mulcaire, P., Ning, Q., Singh, S., Smith, N. A., Subramanian, S., Tsarfaty, R., Wallace, E., Zhang, A., and Zhou, B. Evaluating models’ local decision boundaries via contrast sets, 2020. URL <https://arxiv.org/abs/2004.02709>.
- Gerstgrasser, M., Schaeffer, R., Dey, A., Rafailov, R., Sleight, H., Hughes, J., Korbak, T., Agrawal, R., Pai, D., Gromov, A., Roberts, D. A., Yang, D., Donoho, D. L., and Koyejo, S. Is model collapse inevitable? breaking the curse of recursion by accumulating real and synthetic data, 2024. URL <https://arxiv.org/abs/2404.01413>.
- Haarnoja, T., Zhou, A., Hartikainen, K., Tucker, G., Ha, S., Tan, J., Kumar, V., Zhu, H., Gupta, A., Abbeel, P., and Levine, S. Soft actor-critic algorithms and applications, 2019. URL <https://arxiv.org/abs/1812.05905>.
- Holtzman, A., Buys, J., Du, L., Forbes, M., and Choi, Y. The curious case of neural text degeneration, 2020. URL <https://arxiv.org/abs/1904.09751>.
- Jain, N., Han, K., Gu, A., Li, W.-D., Yan, F., Zhang, T., Wang, Y., Sen, K., Stoica, I., and Gonzalez, J. E. Livecodebench: Holistic and contamination-free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- Katz, M., Kokel, H., and Sreedharan, S. Benchmarking llms on the game of countdown, 2025. URL <https://arxiv.org/abs/2508.02900>.
- Kirk, R., Mediratta, I., Nalmpantis, C., Luketina, J., Hambro, E., Grefenstette, E., and Raileanu, R. Understanding the effects of rlhf on llm generalisation and diversity, 2024. URL <https://arxiv.org/abs/2310.06452>.
- Lanham, T., Chen, A., Radhakrishnan, A., Steiner, B., Denison, C., Hernandez, D., Li, D., Durmus, E., Hubinger, E., Kernion, J., Lukošiūtė, K., Nguyen, K., Cheng, N., Joseph, N., Schiefer, N., Rausch, O., Larson, R., McCandlish, S., Kundu, S., Kadavath, S., Yang, S., Henighan, T., Maxwell, T., Telleen-Lawton, T., Hume, T., Hatfield-Dodds, Z., Kaplan, J., Brauner, J., Bowman, S. R., and Perez, E. Measuring faithfulness in chain-of-thought reasoning, 2023. URL <https://arxiv.org/abs/2307.13702>.
- Li, H. and Klabjan, D. Reverse prompt engineering, 2025. URL <https://arxiv.org/abs/2411.06729>.
- Li, J., Galley, M., Brockett, C., Gao, J., and Dolan, B. A diversity-promoting objective function for neural conversation models, 2016. URL <https://arxiv.org/abs/1510.03055>.

- 495 Li, R., Allal, L. B., Zi, Y., Muennighoff, N., Kocetkov,
496 D., Mou, C., Marone, M., Akiki, C., Li, J., Chim, J.,
497 et al. Taco: Topics in algorithmic code generation. *arXiv*
498 *preprint arXiv:2312.14852*, 2023.
- 499 Liu, L., Wang, Z., Li, L., Xu, C., Lu, Y., Liu, H., Sil, A.,
500 and Li, M. Unary feedback as observation: Incentivizing
501 self-reflection in large language models via multi-turn RL,
502 2026. URL [https://openreview.net/forum?](https://openreview.net/forum?id=GofFHWumFW)
503 [id=GofFHWumFW](https://openreview.net/forum?id=GofFHWumFW).
- 504 Liu, Z., Chen, C., Li, W., Qi, P., Pang, T., Du, C., Lee,
505 W. S., and Lin, M. Understanding rl-zero-like training:
506 A critical perspective, 2025. URL [https://arxiv.](https://arxiv.org/abs/2503.20783)
507 [org/abs/2503.20783](https://arxiv.org/abs/2503.20783).
- 508 Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L.,
509 Wiegrefe, S., Alon, U., Dziri, N., Prabhunoye, S., Yang,
510 Y., Gupta, S., Majumder, B. P., Hermann, K., Welleck,
511 S., Yazdanbakhsh, A., and Clark, P. Self-refine: Iterative
512 refinement with self-feedback, 2023. URL [https://](https://arxiv.org/abs/2303.17651)
513 arxiv.org/abs/2303.17651.
- 514 Mattern, J., Jaghouar, S., Basra, M., Straube, J., Di Fer-
515 rante, M., Gabriel, F., Ong, J. M., Weisser, V., and
516 Hagemann, J. Synthetic-1: Two million collab-
517 oratively generated reasoning traces from deepseek-
518 rl. [https://www.primeintellect.ai/blog/](https://www.primeintellect.ai/blog/synthetic-1-release)
519 [synthetic-1-release](https://www.primeintellect.ai/blog/synthetic-1-release), 2025. Prime Intellect
520 dataset release.
- 521 Meta Llama. Llama 3.2 3b model card, 2024.
522 URL [https://huggingface.co/meta-llama/](https://huggingface.co/meta-llama/Llama-3.2-3B)
523 [Llama-3.2-3B](https://huggingface.co/meta-llama/Llama-3.2-3B). Accessed 2026-01-28.
- 524 Montahaei, E., Alihosseini, D., and Baghshah, M. S. Jointly
525 measuring diversity and quality in text generation mod-
526 els, 2019. URL [https://arxiv.org/abs/1904.](https://arxiv.org/abs/1904.03971)
527 [03971](https://arxiv.org/abs/1904.03971).
- 528 Morris, J. X., Zhao, W., Chiu, J. T., Shmatikov, V., and
529 Rush, A. M. Language model inversion, 2023. URL <https://arxiv.org/abs/2311.13647>.
- 530 Moskovitz, T., Singh, A. K., Strouse, D., Sandholm, T.,
531 Salakhutdinov, R., Dragan, A. D., and McAleer, S. Con-
532 fronting reward model overoptimization with constrained
533 rlhf, 2023. URL [https://arxiv.org/abs/2310.](https://arxiv.org/abs/2310.04373)
534 [04373](https://arxiv.org/abs/2310.04373).
- 535 Nakano, R., Hilton, J., Balaji, S., Wu, J., Ouyang, L.,
536 Kim, C., Hesse, C., Jain, S., Kosaraju, V., Saunders,
537 W., Jiang, X., Cobbe, K., Eloundou, T., Krueger, G.,
538 Button, K., Knight, M., Chess, B., and Schulman, J. Web-
539 gpt: Browser-assisted question-answering with human
540 feedback, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2112.09332)
541 [2112.09332](https://arxiv.org/abs/2112.09332).
- 542 O’Mahony, L., Grinsztajn, L., Schoelkopf, H., and Bider-
543 man, S. Attributing mode collapse in the fine-tuning of
544 large language models. In *ICLR 2024 Workshop on Math-
545 ematical and Empirical Understanding of Foundation*
546 *Models*, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=3pDMYjpOxk)
547 [forum?id=3pDMYjpOxk](https://openreview.net/forum?id=3pDMYjpOxk).
- 548 Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright,
549 C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K.,
Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L.,
Simens, M., Askell, A., Welinder, P., Christiano, P., Leike,
J., and Lowe, R. Training language models to follow
instructions with human feedback, 2022. URL <https://arxiv.org/abs/2203.02155>.
- Pillutla, K., Swayamdipta, S., Zellers, R., Thackstun, J.,
Welleck, S., Choi, Y., and Harchaoui, Z. Mauve: Mea-
suring the gap between neural text and human text using
divergence frontiers, 2021. URL [https://arxiv.](https://arxiv.org/abs/2102.01454)
[org/abs/2102.01454](https://arxiv.org/abs/2102.01454).
- Qi, P., Liu, Z., Zhou, X., Pang, T., Du, C., Lee, W. S.,
and Lin, M. Defeating the training-inference mismatch
via fp16, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2510.26788)
[2510.26788](https://arxiv.org/abs/2510.26788).
- Qwen Team. Qwen2.5 technical report, 2024. URL <https://arxiv.org/abs/2412.15115>.
- Rafailov, R., Sharma, A., Mitchell, E., Ermon, S., Manning,
C. D., and Finn, C. Direct preference optimization: Your
language model is secretly a reward model, 2024. URL <https://arxiv.org/abs/2305.18290>.
- Ribeiro, M. T., Wu, T., Guestrin, C., and Singh, S. Beyond
accuracy: Behavioral testing of nlp models with check-
list, 2020. URL [https://arxiv.org/abs/2005.](https://arxiv.org/abs/2005.04118)
[04118](https://arxiv.org/abs/2005.04118).
- Romoff, J., Henderson, P., Piché, A., Francois-Lavet, V.,
and Pineau, J. Reward estimation for variance reduction
in deep reinforcement learning, 2018. URL <https://arxiv.org/abs/1805.03359>.
- Schrader, M.-P. B. gym-sokoban, 2018. URL [https://](https://github.com/mpSchrader/gym-sokoban)
github.com/mpSchrader/gym-sokoban. Ac-
cessed 2026-01-29.
- Schulman, J., Levine, S., Moritz, P., Jordan, M. I., and
Abbeel, P. Trust region policy optimization, 2017a. URL <https://arxiv.org/abs/1502.05477>.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and
Klimov, O. Proximal policy optimization algorithms,
2017b. URL [https://arxiv.org/abs/1707.](https://arxiv.org/abs/1707.06347)
[06347](https://arxiv.org/abs/1707.06347).

- Schulman, J., Moritz, P., Levine, S., Jordan, M., and Abbeel, P. High-dimensional continuous control using generalized advantage estimation, 2018. URL <https://arxiv.org/abs/1506.02438>.
- Semeniuta, S., Severyn, A., and Gelly, S. On accurate evaluation of gans for language generation, 2019. URL <https://arxiv.org/abs/1806.04936>.
- Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y. K., Wu, Y., and Guo, D. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL <https://arxiv.org/abs/2402.03300>.
- Sheng, G., Zhang, C., Ye, Z., Wu, X., Zhang, W., Zhang, R., Peng, Y., Lin, H., and Wu, C. Hybridflow: A flexible and efficient rlhf framework, 2024. URL <https://arxiv.org/abs/2409.19256>.
- Shinn, N., Cassano, F., Berman, E., Gopinath, A., Narasimhan, K., and Yao, S. Reflexion: Language agents with verbal reinforcement learning, 2023. URL <https://arxiv.org/abs/2303.11366>.
- Shumailov, I., Shumaylov, Z., Zhao, Y., Papernot, N., Anderson, R., and Gal, Y. Ai models collapse when trained on recursively generated data. *Nature*, 631:755–759, 2024. doi: 10.1038/s41586-024-07566-y. URL <https://doi.org/10.1038/s41586-024-07566-y>.
- Siegel, N. Y., Camburu, O.-M., Heess, N., and Perez-Ortiz, M. The probabilities also matter: A more faithful metric for faithfulness of free-text explanations in large language models, 2024. URL <https://arxiv.org/abs/2404.03189>.
- Stiennon, N., Ouyang, L., Wu, J., Ziegler, D. M., Lowe, R., Voss, C., Radford, A., Amodei, D., and Christiano, P. Learning to summarize from human feedback, 2022. URL <https://arxiv.org/abs/2009.01325>.
- Sun, H., Haider, M., Zhang, R., Yang, H., Qiu, J., Yin, M., Wang, M., Bartlett, P., and Zanette, A. Fast best-of-n decoding via speculative rejection, 2024. URL <https://arxiv.org/abs/2410.20290>.
- Tan, S., Luo, M., Cai, C., Venkat, T., Montgomery, K., Hao, A., Wu, T., Balyan, A., Roongta, M., Wang, C., Li, L. E., Popa, R. A., and Stoica, I. rllm: A framework for post-training language agents. <https://pretty-radio-b75.notion.site/rLLM-A-Framework-for-Post-Training-Language-Agents>, 2025. Notion Blog.
- Tao, L., Kulikov, I., Saha, S., Wang, T., Xu, J., Li, S., Weston, J. E., and Yu, P. Hybrid reinforcement: When reward is sparse, it’s better to be dense, 2025. URL <https://arxiv.org/abs/2510.07242>.
- Tevet, G. and Berant, J. Evaluating the evaluation of diversity in natural language generation, 2021. URL <https://arxiv.org/abs/2004.02990>.
- Turpin, M., Michael, J., Perez, E., and Bowman, S. R. Language models don’t always say what they think: Unfaithful explanations in chain-of-thought prompting, 2023. URL <https://arxiv.org/abs/2305.04388>.
- Uesato, J., Kushman, N., Kumar, R., Song, F., Siegel, N., Wang, L., Creswell, A., Irving, G., and Higgins, I. Solving math word problems with process- and outcome-based feedback, 2022. URL <https://arxiv.org/abs/2211.14275>.
- Wang, G., Xie, Y., Jiang, Y., Mandlekar, A., Xiao, C., Zhu, Y., Fan, L., and Anandkumar, A. Voyager: An open-ended embodied agent with large language models, 2023. URL <https://arxiv.org/abs/2305.16291>.
- Wang, J., Liu, J., Fu, Y., Li, Y., Wang, X., Lin, Y., Yue, Y., Zhang, L., Wang, Y., and Wang, K. Harnessing uncertainty: Entropy-modulated policy gradients for long-horizon llm agents, 2025a. URL <https://arxiv.org/abs/2509.09265>.
- Wang, K., Zhang, P., Wang, Z., Gao, Y., Li, L., Wang, Q., Chen, H., Wan, C., Lu, Y., Yang, Z., Wang, L., Krishna, R., Wu, J., Li, F.-F., Choi, Y., and Li, M. VAGEN: Reinforcing world model reasoning for multi-turn VLM agents. *arXiv preprint arXiv:2510.16907*, 2025b.
- Wang, R. and Ammanabrolu, P. A practitioner’s guide to multi-turn agentic reinforcement learning, 2025. URL <https://arxiv.org/abs/2510.01132>.
- Wang, Z., Wang, K., Wang, Q., Zhang, P., Li, L., Yang, Z., Jin, X., Yu, K., Nguyen, M. N., Liu, L., Gottlieb, E., Lu, Y., Cho, K., Wu, J., Fei-Fei, L., Wang, L., Choi, Y., and Li, M. Ragen: Understanding self-evolution in llm agents via multi-turn reinforcement learning, 2025c. URL <https://arxiv.org/abs/2504.20073>.
- Wei, T., Yang, Y., Xing, J., Shi, Y., Lu, Z., and Ye, D. Gtr: Guided thought reinforcement prevents thought collapse in rl-based vlm agent training, 2025. URL <https://arxiv.org/abs/2503.08525>.
- Xu, W., Zhao, W., Wang, Z., Li, Y.-J., Jin, C., Jin, M., Mei, K., Wan, K., and Metaxas, D. N. Epo: Entropy-regularized policy optimization for llm agents reinforcement learning, 2025. URL <https://arxiv.org/abs/2509.22576>.
- Yao, J., Cheng, R., Wu, X., Wu, J., and Tan, K. C. Diversity-aware policy optimization for large language model reasoning, 2025. URL <https://arxiv.org/abs/2505.23433>.

- 605 Yao, S., Chen, H., Yang, J., and Narasimhan, K. Web-
606 shop: Towards scalable real-world web interaction with
607 grounded language agents. In *Advances in Neural Informa-*
608 *tion Processing Systems*, volume 35, pp. 20744–20757,
609 2022.
- 610 Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan,
611 K., and Cao, Y. React: Synergizing reasoning and acting
612 in language models, 2023. URL <https://arxiv.org/abs/2210.03629>.
- 613
614
615 Yin, B., Wang, Q., Zhang, P., Zhang, J., Wang, K., Wang, Z.,
616 Zhang, J., Chandrasegaran, K., Liu, H., Krishna, R., Xie,
617 S., Li, M., Wu, J., and Fei-Fei, L. Spatial mental modeling
618 from limited views. *arXiv preprint arXiv:2506.21458*,
619 2025.
- 620
621 Yu, L., Jiang, W., Shi, H., Yu, J., Liu, Z., Zhang, Y., Kwok,
622 J. T., Li, Z., Weller, A., and Liu, W. Metamath: Bootstrap
623 your own mathematical questions for large language mod-
624 els, 2023. URL [https://arxiv.org/abs/2309.](https://arxiv.org/abs/2309.12284)
625 [12284](https://arxiv.org/abs/2309.12284).
- 626
627 Yu, Q., Zhang, Z., Zhu, R., Yuan, Y., Zuo, X., Yue, Y.,
628 Dai, W., Fan, T., Liu, G., Liu, L., Liu, X., Lin, H., Lin,
629 Z., Ma, B., Sheng, G., Tong, Y., Zhang, C., Zhang, M.,
630 Zhang, W., Zhu, H., Zhu, J., Chen, J., Chen, J., Wang,
631 C., Yu, H., Song, Y., Wei, X., Zhou, H., Liu, J., Ma, W.-
632 Y., Zhang, Y.-Q., Yan, L., Qiao, M., Wu, Y., and Wang,
633 M. Dapo: An open-source llm reinforcement learning
634 system at scale, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2503.14476)
635 [abs/2503.14476](https://arxiv.org/abs/2503.14476).
- 636
637 Yun, L., An, C., Wang, Z., Peng, L., and Shang, J. The
638 price of format: Diversity collapse in llms, 2025. URL
639 <https://arxiv.org/abs/2505.18949>.
- 640
641 Zaman, K. and Srivastava, S. Is chain-of-thought really
642 not explainability? chain-of-thought can be faithful with-
643 out hint verbalization, 2025. URL <https://arxiv.org/abs/2512.23032>.
- 644
645 Zhang, C., Morris, J. X., and Shmatikov, V. Extracting
646 prompts by inverting llm outputs, 2024. URL <https://arxiv.org/abs/2405.15012>.
- 647
648
649 Zhang, Y., Li, Y., Liu, J., Xu, J., Li, Z., Liu, Q., and Li,
650 H. Beyond precision: Training-inference mismatch is
651 an optimization problem and simple lr scheduling fixes
652 it, 2026. URL [https://arxiv.org/abs/2602.](https://arxiv.org/abs/2602.01826)
653 [01826](https://arxiv.org/abs/2602.01826).
- 654
655 Zhu, K., Zhao, Q., Chen, H., Wang, J., and Xie, X. Prompt-
656 bench: A unified library for evaluation of large language
657 models, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2312.07910)
658 [2312.07910](https://arxiv.org/abs/2312.07910).
- 659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800

A. Extended Related Work

Reasoning collapse and policy degeneracy in closed-loop LM and agent RL training. We study a family of degradation phenomena in closed-loop LLM-agent reinforcement learning that has not yet been uniformly defined, but has been repeatedly reported across settings (DeepSeek AI, 2025; Wei et al., 2025). After the model is updated on self-sampled trajectories over time, it may gradually exhibit *reasoning collapse* and *policy-level degeneracy* (DeepSeek AI, 2025; Wei et al., 2025). Here, *reasoning collapse* mainly refers to the rationales, plans, or explanations becoming increasingly templated and less diverse, while their correspondence to the input goal weakens (Wei et al., 2025; Yao et al., 2025; Yun et al., 2025). In contrast, *policy-level degeneracy* refers to behavioral choices concentrating on a small set of easy-to-reproduce action patterns that yield stable scores, with less exploration and less error correction (Feng et al., 2025; Wang & Ammanabrolu, 2025).

This family of phenomena echoes earlier findings in self-training, self-distillation, and iterative fine-tuning on synthetic or model-generated data. When a model repeatedly trains on its own generated distribution, the feedback loop can gradually narrow the effective data distribution, amplify a few high-probability modes, and suppress long-tail behaviors, even when average quality metrics appear stable (Gerstgrasser et al., 2024; Shumailov et al., 2024). In the agent RL setting, closed-loop optimization on on-policy trajectories introduces additional risks, but these risks do not necessarily appear first as an overt failure of the behavioral policy. Instead, a commonly reported pattern is that, even when the agent’s external behavior remains effective or yields stable rewards, language-level reasoning expressions can become concentrated earlier. Plans and explanations may converge to a few reusable narrative skeletons, and their alignment with the specific input goal can weaken (Wei et al., 2025; Xu et al., 2025). In other words, reasoning-level degeneration can decouple from policy-level degeneracy, and in some settings it may precede it (Wang & Ammanabrolu, 2025). In multi-turn interaction, related work also describes several visible signatures of this degradation family, such as within-task convergence across repeated rollouts, cross-task templating where different prompts share the same planning or rhetorical skeleton, and late-stage degeneration where later turns become more mechanical or more conservative (Wang et al., 2025a; Xu et al., 2025).

Evaluating reasoning diversity, input dependence, and reasoning faithfulness. Prior work on evaluating *reasoning diversity* often answers how different the outputs are, but less directly answers whether these differences are *systematically driven by the input goal*, which can blur the interpretation of template-like degeneration under closed-loop training (Tevet & Berant, 2021; Yun et al., 2025). Concretely, common metrics range from lexical measures such as n -gram statistics and self-BLEU (Li et al., 2016; Zhu et al., 2018), to embedding-based dispersion and distributional distances (Pillutla et al., 2021; Tevet & Berant, 2021), as well as token-level uncertainty proxies and multi-sample coverage or consistency analyses (Montahaei et al., 2019; Semeniuta et al., 2019). These metrics primarily capture overall randomness or within-input variability, and they are often less sensitive to whether the reasoning distribution changes coherently *across* inputs (Semeniuta et al., 2019; Tevet & Berant, 2021). Other evaluation protocols rely on model scoring or human preference judgments to compare overall response quality, but they are not designed to isolate input-conditioned reasoning differences, and they may conflate prompt-coupled variation with prompt-agnostic surface diversity, especially when outputs converge to shared formats (Kirk et al., 2024; Yun et al., 2025). This leaves a gap for scalable evaluation of whether reasoning is *diagnostic of the input*, which is particularly salient in multi-turn, stochastic environments where a fixed agent policy can produce diverse yet reusable templates (Wang & Ammanabrolu, 2025). Recent work has started to probe input dependence via behavioral tests and local boundary checks (Gardner et al., 2020; Ribeiro et al., 2020), prompt robustness benchmarks (Zhu et al., 2024), and retrieval-style output–input matching or prompt reconstruction signals (Morris et al., 2023; Gao et al., 2024; Zhang et al., 2024; Li & Klabjan, 2025). However, a unified and scalable treatment tailored to closed-loop agent RL remains limited, even as algorithmic work continues to address long-horizon stability and collapse (Feng et al., 2025; Yao et al., 2025).

A closely related line studies *reasoning faithfulness* (explanation faithfulness), which asks whether a rationale reflects the true basis of a decision rather than a plausible post-hoc story (Lanham et al., 2023; Turpin et al., 2023; Siegel et al., 2024; Zaman & Srivastava, 2025). Our question is related but not equivalent: faithfulness emphasizes whether reasoning causally supports a particular decision, while we focus on a different degeneration risk in closed-loop optimization, namely whether reasoning gradually becomes *less sensitive to the input* and drifts toward reusable templates, even when local explanations remain self-consistent (Kirk et al., 2024). This motivates our decomposition of reasoning diversity into within-input variability and cross-input dependence, and our scalable proxy for the latter through an information-theoretic lens.

Stabilizing multi-turn Agent RL under closed-loop sampling. To improve training stability when aligning LLMs and LLM-based agents, prior work has proposed a broad set of algorithmic and system-level techniques. These include KL

Table 4. MI proxy family.

Type	Proxy	Formula	Notes
Discrete	Retrieval-Acc	$\frac{1}{PG} \sum_{i,k} \mathbf{1}[\arg \max_j \mathbf{L}_{i,k,j} = i]$	Chance level $1/P$ under template collapse
	Recall@ k	$\frac{1}{PG} \sum_{i,k} \mathbf{1}[i \in \text{top-}k_j(\mathbf{L}_{i,k,j})]$	$k \in \{2, 4, 8\}$
Continuous (raw)	MI-Est	$\frac{1}{PG} \sum_{i,k} (\text{matched}_{i,k} - \text{marginal}_{i,k})$	Per-token; approaches 0 under collapse
	MI-Seq-Est	$\frac{1}{PG} \sum_{i,k} (\mathbf{L}_{i,k,i} - \log \frac{1}{P} \sum_j e^{\mathbf{L}_{i,k,j}})$	Per-sequence; no length normalization
Continuous (z-score)	MI-ZScore	$\frac{1}{PG} \sum_{i,k} \frac{\text{matched}_{i,k} - \text{marginal}_{i,k}}{\sigma_{\text{batch}} + \epsilon}$	Normalized by current-batch marginal std
	MI-ZScore-EMA	$\frac{1}{PG} \sum_{i,k} \frac{\text{matched}_{i,k} - \text{marginal}_{i,k}}{\sigma_{\text{EMA}} + \epsilon}$	$\sigma_{\text{EMA}}^{(t)} = \alpha \sigma_{\text{EMA}}^{(t-1)} + (1-\alpha) \sigma_{\text{batch}}^{(t)}$

control or trust-region style constraints, entropy regularization, clipping and normalization in policy-gradient updates, reward shaping and credit assignment, curriculum design, replay or offline-online mixtures, as well as rejection sampling and best-of- N selection (Schulman et al., 2017a;b; 2018; Haarnoja et al., 2019; Stiennon et al., 2022; Ouyang et al., 2022; Rafailov et al., 2024; Sun et al., 2024; Feng et al., 2025; Wang & Ammanabrolu, 2025; Wang et al., 2025a; Xu et al., 2025; Yao et al., 2025). For multi-step agents, researchers have also explored stepwise rewards and intermediate supervision, imitation-to-RL pipelines, and self-correction or reflection signals to support longer-horizon planning and reduce brittle behaviors (Cobbe et al., 2021; Nakano et al., 2022; Uesato et al., 2022; Madaan et al., 2023; Shinn et al., 2023; Wang et al., 2023; Yao et al., 2023; Dou et al., 2025; Wei et al., 2025).

Despite these advances, many stabilization methods are tuned to prevent optimization collapse or to improve overall reward. When the effective learning signal in the closed loop becomes weak or noisy, these methods do not necessarily prevent drift toward prompt-agnostic templates. For example, if most rollouts for the same prompt receive similar rewards regardless of reasoning quality, then the gradient update carries little information about which reasoning path matters (Moskovitz et al., 2023; O’Mahony et al., 2024; Shumailov et al., 2024; Yun et al., 2025). This motivates methods that explicitly manage the balance between task-specific signal and task-agnostic pressure. We adopt a signal-to-noise view of closed-loop updates: we use within-prompt reward variance as a proxy for signal strength, and we filter low-signal samples to maintain an effective SNR, so that exploration and input-conditioned reasoning are less likely to be washed out over long-horizon multi-turn optimization (Romoff et al., 2018; Shao et al., 2024; Tao et al., 2025; Feng et al., 2025; Yao et al., 2025).

B. Additional MI proxies

We list additional MI proxies in Table 6. All variants are derived from in-batch cross-scoring of reasoning traces against prompts, using matched (per-token log-prob under the true prompt) and marginal (per-token log-prob under the uniform prompt mixture) as base quantities. First-turn variants use only the first agent turn; trajectory variants sample across all turns.

C. Detailed Experimental Settings

C.1. Environments and Tasks

We construct a diverse seven-environment testbed to evaluate LLM agents across complementary axes of decision-making complexity, including planning under irreversible dynamics (Sokoban), long-horizon control with non-deterministic transitions (FrozenLake), multi-step symbolic reasoning in mathematics (MetaMathQA, Countdown), multi-turn search and information synthesis (SearchQA), goal-directed web navigation (WebShop), and program synthesis from input-output specifications (DeepCoder). All environments are synthetic and fully controllable, enabling clean analysis of RL learning from scratch without relying on real-world priors.

Sokoban. We use the puzzle Sokoban (Schrader, 2018) to study multi-turn agent interaction with irreversible dynamics. The agent must push boxes to designated target locations within a grid-based warehouse. Unlike standard navigation tasks, Sokoban is characterized by irreversibility: boxes can only be pushed, not pulled, meaning a single misstep can create unsolvable dead-ends where boxes become permanently stuck against walls or corners. This requires the agent to reason ahead and plan multi-step sequences before committing to actions. The reward signal encourages both efficiency and accuracy: +1 for each box successfully placed on a target, -1 for moving a box off a target, +10 upon task completion, and -0.1 per action as a step penalty. We use procedurally generated puzzles with configurable room dimensions and box

counts to ensure diverse training scenarios.

Frozen Lake. This environment of FrozenLake (Brockman et al., 2016) combines long-horizon decision-making with deterministic transitions. The agent navigates a grid of frozen tiles to reach a goal while avoiding holes that terminate the episode. We use the 2% random rate variant of Frozen Lake, where each intended action is executed at a 98% probability. Rewards are sparse: only successful goal-reaching trials receive a reward of +1, with all other outcomes yielding 0. The combination of sparse rewards and long-horizon planning makes this environment challenging for credit assignment.

MetaMathQA. To evaluate mathematical reasoning capabilities, we include MetaMathQA (Yu et al., 2023), a question-answering task drawn from the MetaMathQA dataset. Each episode presents the agent with a mathematical problem requiring multi-step reasoning—ranging from arithmetic and algebra to word problems and geometry. The agent must produce a final answer, and correctness is determined by exact match with the ground truth. To encourage efficient reasoning, we employ a diminishing reward scheme: correct answers on the first attempt receive full reward (1.0), with rewards halving for each subsequent attempt (0.5, 0.25, ...).

Countdown. Inspired by the numbers game from the TV show “Countdown” (Katz et al., 2025), this environment tests compositional arithmetic reasoning. The agent is given a target number and a set of source numbers, and must construct an arithmetic expression using each source number at most once to reach the target exactly. For example, given target 24 and numbers [1, 5, 6, 7], a valid solution is $6 \times (7 - 5 + 1) + 6$. Rewards distinguish between format correctness and solution correctness: full reward (1.0) for correct solutions, partial reward (0.1) for expressions that use the correct numbers but yield incorrect results, and zero for malformed expressions.

DeepCoder. To evaluate agent capabilities in coding environments, we use DeepCoder, a coding benchmark consisting of competitive programming problems. It was used to train DeepSeek-R1-Distill-Qwen-14B with reinforcement learning. The benchmark draws from three resources: PrimeIntellect (Mattern et al., 2025), TACO(Li et al., 2023), and LiveCodeBench v5 (LCBv5) (Jain et al., 2024). In this environment, agents are required to generate a Python function that solves the given programming problem and passes all hidden and public test cases. During training, rewards are assigned based on the number of test cases successfully passed.

SearchQA. To evaluate multi-turn search and question-answering capabilities, we include SearchQA from the RLLM framework (Tan et al., 2025), specifically the Search R1 variant. This environment requires the agent to perform iterative web search and reasoning to answer open-domain questions. The agent must formulate search queries, extract relevant information from retrieved documents, and synthesize answers across multiple interaction turns. Rewards are based on answer correctness and search efficiency, encouraging the agent to balance exploration breadth with reasoning depth.

WebShop. We use WebShop (Yao et al., 2022), an interactive e-commerce environment for evaluating goal-directed multi-turn decision-making. The agent is presented with a shopping instruction (e.g., “find a red shirt under \$30”) and must navigate a simulated online shopping website by issuing search queries, clicking on products, and selecting appropriate items. The environment features a large action space with realistic product catalogs and requires the agent to perform language understanding, attribute matching, and sequential decision-making. Rewards are assigned based on how well the purchased item matches the specified attributes and constraints.

C.2. Training and Evaluation Setup

We conduct our main experiments using Qwen2.5-3B and train with four policy-gradient variants—PPO, DAPO, GRPO, and Dr.GRPO—for up to 400 rollout-update iterations on NVIDIA GPUs using the veRL framework, with early stopping enabled as described below. Each iteration collects $K = 128$ trajectories per environment, organized as $P = 8$ prompt groups with $G = 16$ parallel samples per prompt.

Episode horizons. To match task structure, the interactive environments (Sokoban, Frozen Lake) use up to 5 interaction turns with 2 actions per turn (10 total actions per trajectory). The single-step reasoning tasks (Countdown, MetaMathQA) use 1 turn with 1 action.

Optimization. We use an update batch size of 32 and a per-GPU minibatch size of 4. Policy optimization uses GAE with $(\gamma, \lambda) = (1.0, 1.0)$ and Adam with $(\beta_1, \beta_2) = (0.9, 0.999)$. The actor learning rate is 1×10^{-6} and the critic learning rate is 1×10^{-5} . We apply entropy regularization with coefficient $\beta = 0.001$. For PPO-based methods, we use asymmetric clipping with $\epsilon_{\text{low}} = 0.2$ and $\epsilon_{\text{high}} = 0.28$. We additionally impose a format penalty of -0.1 when the agent fails to output a valid structured response (e.g., missing `<think>` or `<answer>` tags).

Early stopping. We stop training if either (i) reward-variance collapse is detected—the reward variance drops below 10% of the baseline variance (defined as the mean variance over the first 10 training iterations) for 5 consecutive iterations—or (ii) the validation success rate remains below 1% for 5 consecutive evaluation checkpoints.

Filtering ablation. We compare filtered rollouts with $\text{top}_p = 0.9$ (keeping the top 90% of trajectory groups ranked by reward variance) against an unfiltered setting.

Evaluation. We evaluate on a fixed set of 512 validation prompts per environment and decode with temperature $T = 0.5$ using stochastic sampling. We report success rate as the primary metric across all environments.

D. Filtering Ablation Results

We conduct our filtering experiments using Qwen2.5-3B model on Sokoban environment. We summarize the filtering ablation results in Table 5. Each row reports the absolute value of each metric, with the change relative to a section-specific baseline shown in parentheses. Within each block, the first row labeled *baseline* defines the reference point, and all deltas are computed relative to that baseline. We report four metrics: **Task Performance**, defined as the maximum validation success rate attained during training; **MI Proxy**, measured as retrieval accuracy at the training step where task performance peaks; **Entropy**, an estimate of reasoning entropy at the same step; and **Collapse**, a binary indicator of whether validation success ever falls below 0.01 during training.

Sampling Settings. We first study the interaction between filtering and sampling by varying sampling thresholds while holding the reward-variance (RV) filter fixed. Relative to the $\text{top}_p = 1.0$ baseline, reducing top_p or min_p generally improves task performance while reducing entropy, but with heterogeneous effects on MI retention. In contrast, top_k sampling induces a sharper trade-off: MI proxy is often preserved or improved, while gains in task performance are less consistent. These results indicate that filtering behavior is strongly modulated by the sampling regime, even when the underlying filter metric is unchanged.

Filtering Metrics. Next, we fix the sampling scheme and vary the filtering criterion. Switching between RV, entropy-based, entropy-variance, and length-based filters leads to substantial differences in both peak task performance and MI proxy. In particular, SNR-Aware Filtering consistently achieves strong task performance while better preserving MI compared to entropy-based alternatives. Entropy- and length-based filters either suppress MI or fail to prevent collapse, suggesting that reward variance provides a more stable and informative signal for selecting useful rollouts.

Keep Strategy. Finally, we compare *keep-largest* and *keep-smallest* strategies under the same top_k configuration. As expected, retaining high-variance trajectory groups yields substantially higher task performance and MI proxy, while keeping the smallest-variance groups degrades both and markedly increases entropy. This asymmetry supports the hypothesis that high-variance rollouts contain more informative training signal, whereas low-variance rollouts are largely uninformative or noisy.

Summary. Overall, the ablation reveals strong interactions between sampling strategy and filtering choice. More aggressive filtering is not universally beneficial, and the choice of filtering metric is critical: reward-variance filtering consistently improves task performance while maintaining information content, whereas entropy-based heuristics are less reliable and more prone to collapse.

E. Additional Experimental Visualizations

This section presents supplementary visualizations that provide deeper insights into the mechanisms and diagnostics discussed in the main paper. These figures complement the core experimental results with detailed breakdowns of gradient dynamics, diagnostic validity, and reward distribution patterns.

E.1. MI Proxy Metrics During Training

Figure 10 presents six alternative mutual-information proxy metrics tracked over the course of training, complementing the retrieval accuracy shown in Figure 5 of the main paper. All proxies exhibit a consistent pattern: under the *No Filtering* baseline, MI proxies degrade sharply as training progresses, while the three intervention strategies (entropy regularization, KL regularization, and top- p filtering) maintain stable information retention throughout training.

Table 5. Ablation results for sampling strategies, filtering metrics, and keep strategies. Values in parentheses denote the change relative to the corresponding baseline in each block. A crossmark in the Stable column indicates training collapse.

EXPERIMENT SETUP	TASK PERF	MI PROXY	ENTROPY	STABLE
Sampling Strategies				
Top-p = 1.0 (Baseline)	0.17	0.54	2.76	✗
Top-p = 0.9	0.38 (+0.20)	0.84 (+0.29)	1.64 (-1.12)	✓
Top-p = 0.5	0.29 (+0.12)	0.83 (+0.29)	1.88 (-0.88)	✓
Min-p = 0.05	0.42 (+0.25)	0.67 (+0.13)	1.64 (-1.12)	✓
Min-p = 0.2	0.45 (+0.27)	0.36 (-0.18)	3.01 (+0.26)	✓
Top-k = 0.25	0.22 (+0.05)	0.86 (+0.32)	1.28 (-1.48)	✓
Top-k = 0.5	0.44 (+0.27)	0.89 (+0.35)	1.47 (-1.29)	✓
Filtering Metrics				
No Filter (Baseline)	0.17	0.54	2.76	✗
Reward Variance	0.38 (+0.20)	0.84 (+0.29)	1.64 (-1.12)	✓
Reward Sum	0.24 (+0.07)	0.80 (+0.26)	4.18 (+1.42)	✗
Entropy	0.20 (+0.02)	0.41 (-0.14)	2.20 (-0.56)	✗
Entropy Variance	0.23 (+0.06)	0.70 (+0.16)	2.94 (+0.18)	✗
Length	0.16 (-0.02)	0.91 (+0.36)	1.65 (-1.10)	✗
Keep Strategies				
Keep Largest (Baseline)	0.44	0.89	1.47	✓
Keep Smallest	0.29 (-0.15)	0.47 (-0.42)	5.31 (+3.84)	✓

F. Notation and basic identities

F.1. Random variables and distributions

Definition F.1 (Prompts, trajectories, and rollouts). Let X denote an input prompt and Z a reasoning trajectory. A rollout sample is

$$\xi = (x, z, r),$$

with x the prompt, z the realized trajectory, and $r \in \mathbb{R}$ the scalar reward.

We write $\pi_\theta(z | x)$ as the policy and $P(X)$ the prompt distribution. Rollouts are generated by

$$x \sim P(X), \quad z \sim \pi_\theta(\cdot | x), \quad r = R(z; x),$$

where $R(z; x)$ is the reward function.

Definition F.2 (Baseline and advantage). Let $b(x)$ be any function of x only. Define the advantage

$$A(z; x) := R(z; x) - b(x).$$

A standard choice is the conditional-mean baseline $b(x) := \mathbb{E}[R(Z; x) | X = x]$. Then the advantage is zero-mean within each prompt:

$$\mathbb{E}[A(Z; x) | X = x] = \mathbb{E}[R(Z; x) | X = x] - b(x) = 0.$$

Definition F.3 (Score function). Define the score function

$$s(z; x) := \nabla_\theta \log \pi_\theta(z | x).$$

It satisfies the normalization identity

$$\mathbb{E}_{z \sim \pi_\theta(\cdot | x)}[s(z; x)] = \nabla_\theta \int \pi_\theta(z | x) dz = 0.$$

Definition F.4 (Within-prompt reward variance). We quantify within-prompt variation of observed rewards across rollouts by

$$RV(x) := \text{Var}(R(Z; x) | X = x), \quad Z \sim \pi_\theta(\cdot | x).$$

Low $RV(x)$ implies rewards are nearly constant within the prompt, so rollouts are weakly distinguishable by the reward signal. High $RV(x)$ indicates large within-prompt variation of observed rewards which may arise from trajectory-dependent signal or evaluation noise.

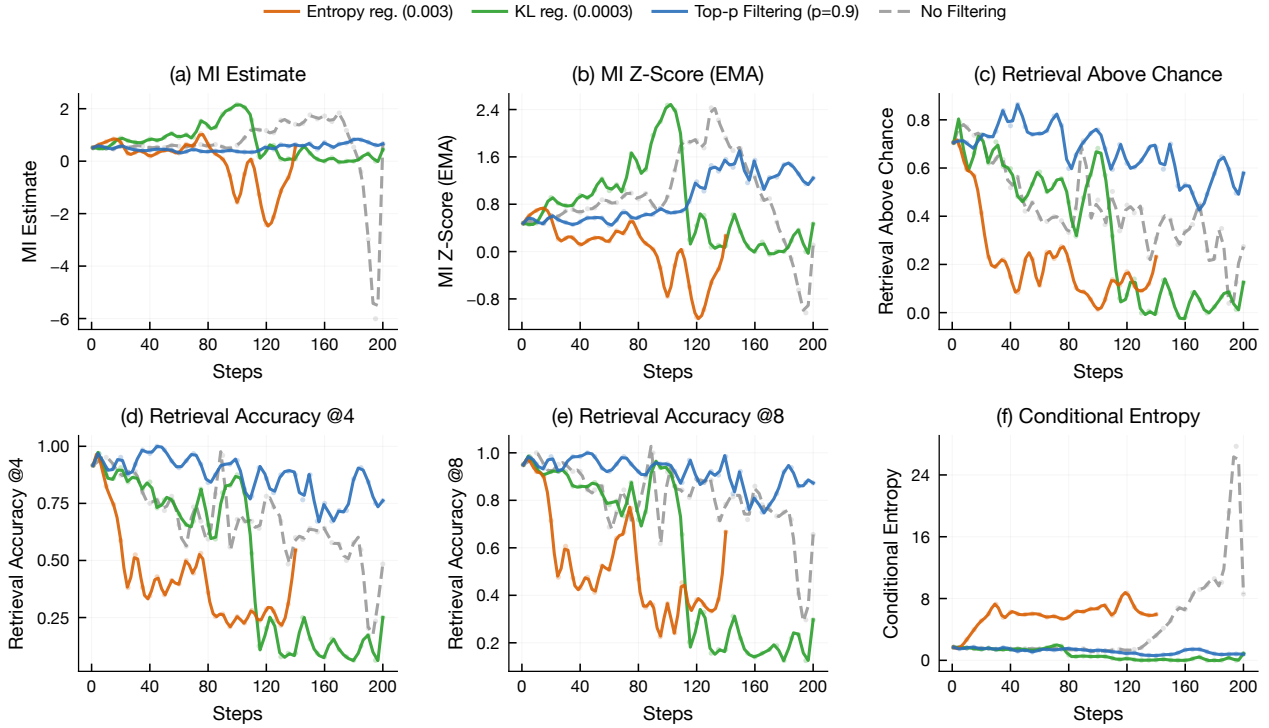


Figure 10. Six MI proxy metrics over training steps. (a) MI Estimate, (b) MI Z-Score (EMA), (c) Retrieval Above Chance, (d) Retrieval Accuracy @4, (e) Retrieval Accuracy @8, (f) Conditional Entropy $H(Z|X)$. Consistent with Figure 5, all proxies confirm that without filtering, MI degrades early, signaling reasoning collapse. Filtering effectively preserves information retention across all metrics, with top- p SNR-aware filtering best maintaining reasoning diversity throughout training.

F.2. Entropy and mutual information

Definition F.5 (Conditional entropy). The within-input variability of reasoning is measured by

$$H(Z | X) := \mathbb{E}_{x \sim P(X)} [H(Z | X = x)] = -\mathbb{E}_{x \sim P(X), z \sim \pi_\theta(\cdot | x)} [\log \pi_\theta(z | x)].$$

The cross-input dependence of reasoning is measured by

$$I(X; Z) := \mathbb{E}_{x \sim P(X), z \sim \pi_\theta(\cdot | x)} \left[\log \frac{\pi_\theta(z | x)}{p_\theta(z)} \right], \quad p_\theta(z) := \mathbb{E}_{x \sim P(X)} [\pi_\theta(z | x)].$$

Equivalently, $I(X; Z) = \mathbb{E}_{x \sim P(X)} [\text{KL}(\pi_\theta(\cdot | x) \| p_\theta)]$.

Decomposition identity (Shannon quantities). For the true distribution induced by π_θ , the Shannon identity

$$H(Z) = H(Z | X) + I(X; Z), \quad (3)$$

serves only as conceptual equation: it specifies the two components we aim to track (within-prompt variability and cross-prompt dependence). In practice we replace these Shannon quantities by scorer-defined proxies, e.g.,

$$\widehat{D}_q := \widehat{\text{NLL}}_q(Z | X) + \widehat{I}_q(X; Z),$$

which is in log-likelihood units under q and does not in general satisfy the Shannon identity unless q matches the evaluated distribution.

Interpretation for reasoning diversity. In our setting, Z is a proxy for a reasoning process (e.g., a chain-of-thought trajectory). A relative decrease in $H(Z | X)$ indicates within-prompt concentration of $\pi_\theta(\cdot | x)$ (entropy collapse). A relative decrease in $I(X; Z)$ indicates weakened input dependence, i.e., trajectories become less diagnostic of x . In our analysis, this can occur when reward-driven updates are weak (e.g., low $\text{RV}(x)$) and the total update is dominated by *reward-agnostic* components (e.g., KL/entropy regularizers). We therefore track these two axes separately; in experiments we use scorer-defined proxies for $H(Z | X)$ and $I(X; Z)$.

G. Scorer-based Proxies for Reasoning Diversity

G.1. Setup and notation

We define scorer-based proxies using a fixed collection of prompts and multiple rollouts per prompt. Throughout this appendix, the scorer q is fixed and used for evaluation.

Definition G.1 (Prompt groups). Using the notation from Definition F.1, sample P prompts $\{x_i\}_{i=1}^P \sim P(X)$. For each prompt x_i , sample G trajectories

$$z_{i,k} \sim \pi_\theta(\cdot | x_i), \quad k = 1, \dots, G.$$

We refer to the set $\{z_{i,k}\}_{k=1}^G$ as a *prompt group*.

Definition G.2 (Teacher-forced scorer and matched-pair score). Let q be a fixed language model used to score how compatible a trajectory z is with a prompt x . Define the matched-pair score

$$\ell_i(z) := \log q(z | x_i).$$

All proxies in this appendix are built from $\ell_i(z)$ and therefore are measured in log-likelihood units under q .

Definition G.3 (Mixture score across prompts). We evaluate each trajectory z under all prompts $\{x_j\}_{j=1}^P$ and define the mixture score

$$\ell_{\text{mix}}(z) := \log \left(\frac{1}{P} \sum_{j=1}^P \exp(\ell_j(z)) \right) = \log \left(\frac{1}{P} \sum_{j=1}^P q(z | x_j) \right).$$

This is the log-likelihood of z under the uniform mixture over prompts induced by q . Equivalently, $\ell_{\text{mix}}(z) = \log(\frac{1}{P} \sum_{j=1}^P q(z | x_j))$ is the log-probability of z under the empirical prompt mixture.

The quantities defined above depend on the sampled prompt set $\{x_i\}_{i=1}^P$ and on the fixed scorer q . They are proxies for within-prompt variability and input dependence of trajectories, and should not be interpreted as exact Shannon entropies or mutual information unless q matches the evaluated conditional distribution.

H. Formal Definition of the Filtering Operator

Definition H.1 (Filtering operator). Let \mathcal{B} be a minibatch of samples. A *filtering operator* is specified by:

(i) **Grouping key.** A grouping function $g : \mathcal{B} \rightarrow \mathcal{G}$ that assigns each sample $\xi \in \mathcal{B}$ a group label

$$u = g(\xi).$$

For $u \in \mathcal{G}$, define the induced group subset

$$\mathcal{B}_u := \{\xi \in \mathcal{B} : g(\xi) = u\}.$$

(ii) **Group statistic.** A statistic $\phi : \mathcal{B} \rightarrow \mathbb{R}$ that depends only on the samples in the group, and we write $\phi(\mathcal{B}_u)$ for the value computed from \mathcal{B}_u .

(iii) **Selection rule (mask).** Given a threshold $\tau \in \mathbb{R}$, the binary mask is

$$m(u) := \mathbf{1}\{\phi(\mathcal{B}_u) \geq \tau\}.$$

(iv) **Filtered objective.** For a per-sample RL loss $L_\theta(\xi)$, the filtered objective is

$$\mathcal{L}_{\text{filt}}(\theta) = \frac{1}{|\mathcal{B}|} \sum_{\xi \in \mathcal{B}} m(g(\xi)) L_\theta(\xi).$$

Remark (post-sampling). Filtering is applied after sampling and only masks gradients; it does not change the rollout distribution.

Remark (normalization). In practice one may normalize by the number of kept samples or kept groups (instead of $|\mathcal{B}|$), which rescales the gradient but does not change which samples contribute nonzero gradients.

H.1. Filtering Strategy Variants

We compare multiple filtering strategies for selecting high-signal prompt groups. All variants share the same grouping structure (prompts with G rollouts each) and statistic (reward variance $\widehat{\text{Var}}(R \mid X = x_i)$ for group i), but differ in the selection rule.

Top-p (nucleus-style) filtering. The main method used in this paper. Given keep rate $\rho \in (0, 1]$, rank prompts by descending reward variance and select the smallest prefix whose cumulative variance mass reaches $\rho \sum_i \widehat{\text{Var}}(R \mid X = x_i)$. Formally, let σ be the permutation such that $\widehat{\text{Var}}(R \mid X = x_{\sigma(1)}) \geq \dots \geq \widehat{\text{Var}}(R \mid X = x_{\sigma(N)})$, and define

$$k^* = \min \left\{ k : \sum_{j=1}^k \widehat{\text{Var}}(R \mid X = x_{\sigma(j)}) \geq \rho \sum_{i=1}^N \widehat{\text{Var}}(R \mid X = x_i) \right\}.$$

The kept set is $S = \{\sigma(1), \dots, \sigma(k^*)\}$. This adaptive selection concentrates updates on high-variance prompts while automatically adjusting the kept count based on the variance distribution. When the batch contains many near-zero-variance prompts, top-p can reject the entire batch if the threshold cannot be reached, providing a natural safeguard against degenerate updates.

Top-k (proportional) filtering. An alternative fixed-proportion baseline. Given $\rho \in (0, 1]$, compute $k = \lfloor \rho N \rfloor$ and select the top k prompts by reward variance:

$$S = \{\sigma(1), \dots, \sigma(k)\}.$$

Unlike top-p, top-k always retains exactly k groups regardless of the variance distribution. This can be less adaptive: when most prompts have near-zero variance, top-k still keeps the highest-variance subset even if all retained prompts carry weak signal.

Min-p (threshold) filtering. Inspired by min-p sampling, this strategy keeps all prompts whose variance exceeds a fraction of the maximum variance. Given threshold parameter $p \in (0, 1]$, define

$$\tau = p \cdot \max_i \widehat{\text{Var}}(R \mid X = x_i),$$

and keep all groups above the threshold:

$$S = \{i : \widehat{\text{Var}}(R \mid X = x_i) \geq \tau\}.$$

This directly enforces a minimum quality bar: only prompts within a factor of p of the best prompt are retained. The kept count varies with the variance distribution, making this method highly adaptive but potentially unstable when the maximum variance fluctuates.

Reverse top-p (low-variance) filtering. A diagnostic baseline that intentionally selects low-variance prompts. Rank prompts by *ascending* reward variance and select the smallest prefix whose cumulative variance mass reaches $\rho \sum_i \widehat{\text{Var}}(R \mid X = x_i)$. This inverted strategy is used in ablation studies to confirm that high variance is essential for effective updates: training on low-variance prompts should degrade both MI and task performance, validating the SNR hypothesis.

Implementation notes. All strategies can be configured to exclude zero-variance groups (setting `include_zero=False`) before selection, which removes prompts where all rollouts received identical rewards. For top-p, we use a small epsilon $\varepsilon = 0.01$ to ensure numerical stability when checking whether the cumulative threshold is reached. Additional implementation details and hyperparameter sensitivity are in the codebase.

I. RV Controls Task-Signal Magnitude and SNR

I.1. Setup

We use the policy/score/baseline/advantage notation from Appendix F.

In particular, for a fixed prompt x we write $z \sim \pi_\theta(\cdot \mid x)$, $s(z; x) = \nabla_\theta \log \pi_\theta(z \mid x)$, $A(z; x) = R(z; x) - b(x)$ with $b(x) = \mathbb{E}[R \mid X = x]$, and $\text{RV}(x) = \text{Var}(R \mid X = x) = \mathbb{E}[A^2 \mid X = x]$.

1100 I.2. Assumption

1101 **Assumption I.1** (Reward decomposition). The observed reward admits a decomposition

$$1102 R(z; x) = \mu(x, z) + \varepsilon, \quad \mu(x, z) := \mathbb{E}[R(z; x) \mid x, z],$$

1103 where $\mu(x, z)$ is the trajectory-dependent mean reward and ε is a zero-mean noise term satisfying

$$1104 \mathbb{E}[\varepsilon \mid x, z] = 0, \quad \text{Var}(\varepsilon \mid x, z) = \sigma^2(x) \geq 0.$$

1105 Moreover, the score $s(z; x) = \nabla_{\theta} \log \pi_{\theta}(z \mid x)$ is a deterministic (measurable) function of (x, z) .

1106 I.3. Task-gradient magnitude is RV-controlled

1107 The next result shows that the task-gradient norm for a given prompt is at most proportional to the square root of its within-prompt reward variance $\text{RV}(x)$. In particular, when $\text{RV}(x)$ is small, the task gradient is provably weak.

1108 **Theorem I.2** (Task gradient magnitude is RV-controlled). *Assume the baseline is the conditional mean $b(x) = \mathbb{E}[R \mid X = x]$, and $g_{\text{task}}(x) := \mathbb{E}[A(z; x) s(z; x) \mid X = x]$. Then*

$$1109 \|g_{\text{task}}(x)\| \leq \sqrt{\text{RV}(x)} \sqrt{\mathbb{E}[\|s(z; x)\|^2 \mid X = x]}.$$

1110 *Proof.* Fix a prompt x and take randomness over $z \sim \pi_{\theta}(\cdot \mid x)$. For brevity write $A := A(z; x)$ and $s := s(z; x)$. Then

$$1111 g_{\text{task}}(x) = \mathbb{E}[A s \mid X = x].$$

1112 For any unit vector $u \in \mathbb{R}^d$ with $\|u\| = 1$,

$$1113 |\langle u, g_{\text{task}}(x) \rangle| = |\mathbb{E}[A \langle u, s \rangle \mid X = x]| \leq \sqrt{\mathbb{E}[A^2 \mid X = x]} \sqrt{\mathbb{E}[\langle u, s \rangle^2 \mid X = x]},$$

1114 where the inequality is Cauchy-Schwarz. Moreover, $\langle u, s \rangle^2 \leq \|u\|^2 \|s\|^2 = \|s\|^2$, hence

$$1115 |\langle u, g_{\text{task}}(x) \rangle| \leq \sqrt{\mathbb{E}[A^2 \mid X = x]} \sqrt{\mathbb{E}[\|s\|^2 \mid X = x]}.$$

1116 Taking the supremum over all unit vectors u yields

$$1117 \|g_{\text{task}}(x)\| \leq \sqrt{\mathbb{E}[A^2 \mid X = x]} \sqrt{\mathbb{E}[\|s\|^2 \mid X = x]}.$$

1118 Finally, with $b(x) = \mathbb{E}[R \mid X = x]$ we have $\mathbb{E}[A \mid X = x] = 0$ and thus

$$1119 \mathbb{E}[A^2 \mid X = x] = \text{Var}(R \mid X = x) = \text{RV}(x).$$

1120 Substituting completes the proof. \square

1121 I.4. SNR is upper bounded by RV and reward noise

1122 The following theorem shows that the signal-to-noise ratio of the G -sample Monte Carlo gradient estimator is upper-bounded by $\sqrt{G} \cdot \sqrt{\text{RV}(x)}/\sigma(x)$. When reward variance is low relative to reward noise, the estimator is dominated by noise.

1123 **Theorem I.3** (SNR upper bound by RV and noise). *Let $\hat{g}_{\text{task}}(x)$ be the G -sample Monte Carlo estimator*

$$1124 \hat{g}_{\text{task}}(x) := \frac{1}{G} \sum_{k=1}^G A_k s_k, \quad A_k := A(z_k; x), \quad s_k := s(z_k; x),$$

1125 with $z_1, \dots, z_G \stackrel{\text{i.i.d.}}{\sim} \pi_{\theta}(\cdot \mid x)$. Define

$$1126 \text{SNR}(x) := \frac{\|g_{\text{task}}(x)\|}{\sqrt{\mathbb{E}[\|\hat{g}_{\text{task}}(x) - g_{\text{task}}(x)\|^2 \mid X = x]}}.$$

1127 Under Assumption I.1 and with baseline $b(x) = \mathbb{E}[R \mid X = x]$,

$$1128 \text{SNR}(x) \leq \sqrt{G} \cdot \frac{\sqrt{\text{RV}(x)}}{\sigma(x)}.$$

1129 If $\sigma(x) = 0$, the bound is vacuous.

1130 *Proof.* Fix a prompt x . Let $z_1, \dots, z_G \stackrel{\text{i.i.d.}}{\sim} \pi_{\theta}(\cdot \mid x)$ and write

$$1131 \hat{g} = \frac{1}{G} \sum_{k=1}^G A_k s_k, \quad g = \mathbb{E}[A s \mid x],$$

where $(A_k, s_k) = (A(z_k; x), s(z_k; x))$ and $(A, s) = (A(z; x), s(z; x))$ for $z \sim \pi_\theta(\cdot | x)$.

Let $Y_k := A_k s_k$. Then $\hat{g} = \frac{1}{G} \sum_{k=1}^G Y_k$ and $g = \mathbb{E}[Y_1 | x]$, hence

$$\hat{g} - g = \frac{1}{G} \sum_{k=1}^G (Y_k - g).$$

Using i.i.d. conditional on x ,

$$\begin{aligned} \mathbb{E}[\|\hat{g} - g\|^2 | x] &= \frac{1}{G^2} \mathbb{E} \left[\left\| \sum_{k=1}^G (Y_k - g) \right\|^2 \middle| x \right] \\ &= \frac{1}{G^2} \sum_{k=1}^G \mathbb{E}[\|Y_k - g\|^2 | x] + \frac{1}{G^2} \sum_{k \neq \ell} \mathbb{E}[\langle Y_k - g, Y_\ell - g \rangle | x] \\ &= \frac{1}{G^2} \sum_{k=1}^G \mathbb{E}[\|Y_k - g\|^2 | x] \\ &= \frac{1}{G} \mathbb{E}[\|As - g\|^2 | x]. \end{aligned}$$

Under Assumption I.1 and with baseline $b(x) = \mathbb{E}[R | X = x]$, write $R = \mu + \varepsilon$ with $\mu(x, z) = \mathbb{E}[R | x, z]$. Since $b(x) = \mathbb{E}[R | x] = \mathbb{E}[\mu | x]$,

$$A = R - b(x) = (\mu - \mathbb{E}[\mu | x]) + \varepsilon =: A_\mu + \varepsilon.$$

Using $A = A_\mu + \varepsilon$,

$$As - g = (A_\mu s - g) + \varepsilon s,$$

so

$$\|As - g\|^2 = \|A_\mu s - g\|^2 + \|\varepsilon s\|^2 + 2\langle A_\mu s - g, \varepsilon s \rangle.$$

Moreover,

$$\mathbb{E}[\langle A_\mu s - g, \varepsilon s \rangle | x] = \mathbb{E} \left[\mathbb{E}[\langle A_\mu s - g, \varepsilon s \rangle | x, z] \middle| x \right] = \mathbb{E}[\langle A_\mu s - g, s \rangle \mathbb{E}[\varepsilon | x, z] | x] = 0,$$

hence

$$\mathbb{E}[\|As - g\|^2 | x] \geq \mathbb{E}[\|\varepsilon s\|^2 | x].$$

Combining with the variance decomposition above,

$$\mathbb{E}[\|\hat{g} - g\|^2 | x] \geq \frac{1}{G} \mathbb{E}[\|\varepsilon s\|^2 | x].$$

Since $\|\varepsilon s\|^2 = \varepsilon^2 \|s\|^2$ and s is measurable given (x, z) ,

$$\begin{aligned} \mathbb{E}[\|\varepsilon s\|^2 | x] &= \mathbb{E} \left[\mathbb{E}[\varepsilon^2 \|s\|^2 | x, z] \middle| x \right] \\ &= \mathbb{E} \left[\|s\|^2 \mathbb{E}[\varepsilon^2 | x, z] \middle| x \right] \\ &= \mathbb{E} \left[\|s\|^2 \sigma^2(x) \middle| x \right] = \sigma^2(x) \mathbb{E}[\|s\|^2 | x], \end{aligned}$$

where $\mathbb{E}[\varepsilon^2 | x, z] = \text{Var}(\varepsilon | x, z) = \sigma^2(x)$ by Assumption I.1. Therefore

$$\mathbb{E}[\|\hat{g} - g\|^2 | x] \geq \frac{1}{G} \sigma^2(x) \mathbb{E}[\|s\|^2 | x].$$

By Theorem I.2,

$$\|g\| = \|\mathbb{E}[As | x]\| \leq \sqrt{\text{RV}(x)} \sqrt{\mathbb{E}[\|s\|^2 | x]}.$$

Thus, with $\text{SNR}(x) := \frac{\|g\|}{\sqrt{\mathbb{E}[\|\hat{g} - g\|^2 | x]}}$,

$$\text{SNR}(x) \leq \frac{\sqrt{\text{RV}(x)} \sqrt{\mathbb{E}[\|s\|^2 | x]}}{\sqrt{\frac{1}{G} \sigma^2(x) \mathbb{E}[\|s\|^2 | x]}} = \sqrt{G} \cdot \frac{\sqrt{\text{RV}(x)}}{\sigma(x)}. \quad \square$$

1210 I.5. Low-SNR updates induce parameter drift

1211 When updates carry no directional signal (zero mean), the parameter drifts away from initialization at a rate linear in the
 1212 number of steps. This illustrates why sustained low-SNR updates are harmful even if they do not systematically push in a
 1213 wrong direction.
 1214

1215 **Theorem I.4** (Illustrative random-walk drift under zero-mean noise). *Consider SGD-style updates*

$$1216 \theta_{t+1} = \theta_t + \eta \xi_t,$$

1217 where $\{\xi_t\}_{t \geq 0}$ are independent, $\mathbb{E}[\xi_t] = 0$, and $\mathbb{E}[\|\xi_t\|^2] = v < \infty$ for all t . Then for any $T \geq 1$,

$$1218 \mathbb{E}[\|\theta_T - \theta_0\|^2] = \eta^2 T v.$$

1219
 1220 *Proof.* Unrolling the recursion yields

$$1221 \theta_T - \theta_0 = \eta \sum_{t=0}^{T-1} \xi_t.$$

1222 Therefore,

$$1223 \|\theta_T - \theta_0\|^2 = \eta^2 \left\| \sum_{t=0}^{T-1} \xi_t \right\|^2 = \eta^2 \left(\sum_{t=0}^{T-1} \|\xi_t\|^2 + 2 \sum_{0 \leq i < j \leq T-1} \langle \xi_i, \xi_j \rangle \right).$$

1224 Taking expectation and using independence with $\mathbb{E}[\xi_t] = 0$,

$$1225 \mathbb{E}\langle \xi_i, \xi_j \rangle = \langle \mathbb{E}[\xi_i], \mathbb{E}[\xi_j] \rangle = 0, \quad i \neq j.$$

1226 Hence the cross terms vanish and

$$1227 \mathbb{E}[\|\theta_T - \theta_0\|^2] = \eta^2 \sum_{t=0}^{T-1} \mathbb{E}[\|\xi_t\|^2] = \eta^2 T v,$$

1228 where we used $\mathbb{E}[\|\xi_t\|^2] = v$ for all t . □

1229 J. Template Mixing Reduces Input Dependence

1230 If the policy's conditional distribution is contaminated by a prompt-independent component $q(z)$ with mixing weight α ,
 1231 the resulting mutual information $I_\alpha(X; Z)$ contracts by at least a factor of $(1 - \alpha)$. This formalizes the intuition that even
 1232 partial drift toward a shared template erodes input dependence.

1233 **Lemma J.1** (Template mixing contracts mutual information). *Let $X \sim P(X)$ and $Z | X = x \sim p(z | x)$ with marginal*
 1234 *$p(z) = \mathbb{E}_{x \sim P}[p(z | x)]$. Fix any prompt-independent distribution $q(z)$. For $\alpha \in [0, 1]$, define the mixed conditional and*
 1235 *marginal*

$$1236 p_\alpha(z | x) := (1 - \alpha)p(z | x) + \alpha q(z), \quad p_\alpha(z) := (1 - \alpha)p(z) + \alpha q(z).$$

1237 Let $I_\alpha(X; Z)$ denote the mutual information under $p_\alpha(x, z) = P(x)p_\alpha(z | x)$. Then

$$1238 I_\alpha(X; Z) \leq (1 - \alpha) I(X; Z).$$

1239 *Proof.* For any fixed x ,

$$1240 \text{KL}(p_\alpha(\cdot | x) \| p_\alpha(\cdot)) = \mathbb{E}_{z \sim p_\alpha(\cdot | x)} \left[\log \frac{p_\alpha(z | x)}{p_\alpha(z)} \right]$$

$$1241 = \mathbb{E}_{z \sim p_\alpha(\cdot | x)} [\log p_\alpha(z | x) - \log p_\alpha(z)].$$

1242 Taking expectation over $x \sim P(x)$ gives

$$1243 I_\alpha(X; Z) = \mathbb{E}_x [\text{KL}(p_\alpha(\cdot | x) \| p_\alpha(\cdot))].$$

1244 The same identity holds for $I(X; Z)$ with p_α replaced by p .

1245 By joint convexity of $\text{KL}(\cdot \| \cdot)$ (Cover & Thomas, 2006, Theorem 2.7.2), for any distributions a, b, c, d and any $\alpha \in [0, 1]$,

$$1246 \text{KL}((1 - \alpha)a + \alpha b \| (1 - \alpha)c + \alpha d) \leq (1 - \alpha)\text{KL}(a \| c) + \alpha \text{KL}(b \| d).$$

1247 Let $a = p(\cdot | x)$, $b = q$, $c = p(\cdot)$, and $d = q$. Since

$$1248 p_\alpha(\cdot | x) = (1 - \alpha)p(\cdot | x) + \alpha q, \quad p_\alpha(\cdot) = (1 - \alpha)p(\cdot) + \alpha q,$$

1265 we obtain

$$1266 \quad \text{KL}(p_\alpha(\cdot | x) \| p_\alpha(\cdot)) \leq (1 - \alpha)\text{KL}(p(\cdot | x) \| p(\cdot)) + \alpha \text{KL}(q \| q)$$

$$1267 \quad = (1 - \alpha)\text{KL}(p(\cdot | x) \| p(\cdot)).$$

1268 Averaging over $x \sim P(x)$ yields

$$1270 \quad \mathbb{E}_x \left[\text{KL}(p_\alpha(\cdot | x) \| p_\alpha(\cdot)) \right] \leq (1 - \alpha) \mathbb{E}_x \left[\text{KL}(p(\cdot | x) \| p(\cdot)) \right].$$

1271 Using the identity $I(X; Z) = \mathbb{E}_x \left[\text{KL}(p(\cdot | x) \| p(\cdot)) \right]$ (and the analogous one for I_α), we obtain

$$1272 \quad I_\alpha(X; Z) \leq (1 - \alpha) I(X; Z),$$

1273 which proves the lemma. □

1274 *Remark J.2.* The continuity bound $f(\varepsilon)$ depends on $\log(|\mathcal{X}||\mathcal{Z}|)$, which can be extremely large for LLM token spaces.
 1275 Therefore, this result should be understood as a qualitative guarantee that KL-closeness implies MI-closeness in principle,
 1276 rather than a tight quantitative bound in practice.

1281

1282 K. Filtering Reduces Gradient-Estimation MSE

1283 K.1. Setup

1284 Consider P groups indexed by $i \in \{1, \dots, P\}$. Group i contains G rollouts, and $\hat{g}_i \in \mathbb{R}^d$ denotes the *group-level* gradient
 1285 estimator (already averaged over the G rollouts in the group). We model

$$1286 \quad \hat{g}_i = g_i + \varepsilon_i, \quad \mathbb{E}[\varepsilon_i] = 0, \quad \mathbb{E}\|\varepsilon_i\|^2 = \sigma_i^2,$$

1287 where $\{\varepsilon_i\}_{i=1}^P$ are independent across groups. For a kept set S of groups, we write $n := |S|$ for the number of kept groups.

1291 K.2. Unfiltered vs. filtered estimators

1292 Define the unfiltered batch estimator and its mean:

$$1293 \quad \hat{\bar{g}} := \frac{1}{P} \sum_{i=1}^P \hat{g}_i, \quad \bar{g} := \frac{1}{P} \sum_{i=1}^P g_i.$$

1294 Let $S \subseteq \{1, \dots, P\}$ be the set of kept groups with $|S| = n$. Define the filtered estimator and its mean:

$$1295 \quad \hat{g}_S := \frac{1}{n} \sum_{i \in S} \hat{g}_i, \quad \bar{g}_S := \frac{1}{n} \sum_{i \in S} g_i.$$

1296 By retaining only a subset of prompt groups, the filtered estimator's mean-squared error depends solely on the noise
 1297 variances of the kept groups. Dropping high-noise (low-RV) groups directly lowers the estimation error.

1300 **Theorem K.1** (MSE of the filtered estimator). \hat{g}_S is unbiased for \bar{g}_S and satisfies

$$1301 \quad \mathbb{E}\|\hat{g}_S - \bar{g}_S\|^2 = \frac{1}{n^2} \sum_{i \in S} \sigma_i^2.$$

1302 *Proof.* By the setup, $\hat{g}_i = g_i + \varepsilon_i$ with $\mathbb{E}[\varepsilon_i] = 0$, hence

$$1303 \quad \mathbb{E}[\hat{g}_i] = g_i.$$

1304 Therefore,

$$1305 \quad \mathbb{E}[\hat{g}_S] = \frac{1}{n} \sum_{i \in S} \mathbb{E}[\hat{g}_i] = \frac{1}{n} \sum_{i \in S} g_i = \bar{g}_S.$$

1306 Moreover,

$$1307 \quad \hat{g}_S - \bar{g}_S = \frac{1}{n} \sum_{i \in S} (\hat{g}_i - g_i) = \frac{1}{n} \sum_{i \in S} \varepsilon_i.$$

1319

Therefore,

$$\begin{aligned} \mathbb{E}\|\widehat{g}_S - \bar{g}_S\|^2 &= \frac{1}{n^2} \mathbb{E}\left\|\sum_{i \in S} \varepsilon_i\right\|^2 \\ &= \frac{1}{n^2} \left(\sum_{i \in S} \mathbb{E}\|\varepsilon_i\|^2 + \sum_{\substack{i, j \in S \\ i \neq j}} \mathbb{E}\langle \varepsilon_i, \varepsilon_j \rangle \right). \end{aligned}$$

By independence and $\mathbb{E}[\varepsilon_i] = 0$, for $i \neq j$ we have

$$\mathbb{E}\langle \varepsilon_i, \varepsilon_j \rangle = \langle \mathbb{E}[\varepsilon_i], \mathbb{E}[\varepsilon_j] \rangle = 0,$$

so the cross terms vanish. Hence

$$\mathbb{E}\|\widehat{g}_S - \bar{g}_S\|^2 = \frac{1}{n^2} \sum_{i \in S} \mathbb{E}\|\varepsilon_i\|^2 = \frac{1}{n^2} \sum_{i \in S} \sigma_i^2.$$

□

Remark (bias relative to the original objective). While \widehat{g}_S is unbiased for the *filtered* mean gradient \bar{g}_S , it is generally biased for the *unfiltered* mean gradient \bar{g} unless S is chosen independently of $\{g_i\}$ or g_i is constant across groups.

L. Reward-Agnostic Regularizers and Update Dominance

L.1. Setup

Similarly, fix a prompt x and consider trajectories $z \sim \pi_\theta(\cdot | x)$ with reward $R(z; x)$ and baseline $b(x)$. Define the reward-driven (task) gradient

$$g_{\text{task}}(x) := \mathbb{E}[(R(z; x) - b(x)) s(z; x) | X = x], \quad s(z; x) := \nabla_\theta \log \pi_\theta(z | x).$$

Let $g_{\text{reg}}(x)$ denote an update component that is computed without multiplying the reward (or advantage), e.g.,

$$g_{\text{reg}}(x) := \lambda_{\text{KL}} g_{\text{KL}}(x) + \lambda_{\text{ent}} g_{\text{ent}}(x),$$

where $g_{\text{KL}}(x)$ and $g_{\text{ent}}(x)$ are gradients of prompt-level distributional regularizers. We write the total expected update as

$$g_{\text{total}}(x) = g_{\text{task}}(x) + g_{\text{reg}}(x).$$

To summarize relative influence, define the dominance ratio

$$\rho(x) := \frac{\|g_{\text{reg}}(x)\|}{\|g_{\text{task}}(x)\| + \|g_{\text{reg}}(x)\|} \in [0, 1].$$

We refer to $g_{\text{reg}}(x)$ as *reward-agnostic* since it does not use within-prompt reward differences to weight trajectories.

L.2. Low-RV prompts amplify regularizer influence

When reward variance is small, the task gradient weakens (by Theorem I.2) while regularizer gradients remain largely flat across prompts. Consequently, the regularizer’s share of the total update grows on low-RV prompts, formalizing why these prompts are more prone to input-agnostic drift.

By Theorem I.2, for any prompt x ,

$$\|g_{\text{task}}(x)\| \leq \sqrt{\text{RV}(x)} \sqrt{\mathbb{E}[\|s\|^2 | X = x]}.$$

Therefore the dominance ratio

$$\rho(x) = \frac{\|g_{\text{reg}}(x)\|}{\|g_{\text{task}}(x)\| + \|g_{\text{reg}}(x)\|}$$

admits the lower bound

$$\rho(x) \geq \frac{\|g_{\text{reg}}(x)\|}{\|g_{\text{reg}}(x)\| + \sqrt{\text{RV}(x)} \sqrt{\mathbb{E}[\|s\|^2 | X = x]}}.$$

In particular, if $\|g_{\text{reg}}(x)\|$ and $\mathbb{E}[\|s\|^2 | X = x]$ vary slowly across prompts compared to $\text{RV}(x)$, then smaller $\text{RV}(x)$ implies larger $\rho(x)$, i.e., the total update is more strongly shaped by reward-agnostic regularizers on low-RV prompts.

M. KL-Closeness to the Base Implies MI-Closeness

If the current policy stays uniformly close to a reference policy in KL divergence, then the mutual information $I(X; Z)$ between inputs and reasoning also remains close. This means strong KL constraints preserve—but do not necessarily increase—input dependence.

Theorem M.1. *To avoid measure-theoretic issues, assume X is supported on a finite set \mathcal{X} and Z takes values in a finite set \mathcal{Z} . Let $P(X)$ be the prompt distribution and define*

$$P_\theta(x, z) := P(x)\pi_\theta(z | x), \quad P_0(x, z) := P(x)\pi_0(z | x).$$

If

$$\sup_{x \in \mathcal{X}} \text{KL}(\pi_\theta(\cdot | x) \| \pi_0(\cdot | x)) \leq \varepsilon,$$

then there exists $f(\varepsilon) \rightarrow 0$ as $\varepsilon \rightarrow 0$ such that

$$|I_\theta(X; Z) - I_0(X; Z)| \leq f(\varepsilon).$$

Proof. By the chain rule for KL divergence,

$$\text{KL}(P_\theta(X, Z) \| P_0(X, Z)) = \mathbb{E}_{x \sim P}[\text{KL}(\pi_\theta(\cdot | x) \| \pi_0(\cdot | x))].$$

Under the assumption $\sup_{x \in \mathcal{X}} \text{KL}(\pi_\theta(\cdot | x) \| \pi_0(\cdot | x)) \leq \varepsilon$, we obtain

$$\text{KL}(P_\theta(X, Z) \| P_0(X, Z)) \leq \varepsilon.$$

By Pinsker’s inequality,

$$\|P_\theta(X, Z) - P_0(X, Z)\|_{\text{TV}} \leq \sqrt{\frac{1}{2} \text{KL}(P_\theta(X, Z) \| P_0(X, Z))} \leq \sqrt{\frac{\varepsilon}{2}} =: \delta.$$

Since $\|P_\theta(X, Z) - P_0(X, Z)\|_{\text{TV}} \leq \delta$ and (X, Z) takes values in a finite alphabet $\mathcal{X} \times \mathcal{Z}$, the Fannes-Audenaert inequality implies

$$|H_\theta(X, Z) - H_0(X, Z)| \leq \delta \log(|\mathcal{X}||\mathcal{Z}| - 1) + h_2(\delta),$$

where $H_\theta(\cdot)$ denotes entropy under P_θ , and $h_2(\cdot)$ is the binary entropy. Moreover, total variation does not increase under marginalization, so

$$\|P_\theta(Z) - P_0(Z)\|_{\text{TV}} \leq \delta,$$

and applying Fannes-Audenaert on the alphabet \mathcal{Z} yields

$$|H_\theta(Z) - H_0(Z)| \leq \delta \log(|\mathcal{Z}| - 1) + h_2(\delta) \leq \delta \log(|\mathcal{X}||\mathcal{Z}| - 1) + h_2(\delta).$$

Finally, using $I(X; Z) = H(X) + H(Z) - H(X, Z)$ and noting that $P_\theta(X) = P_0(X) = P(X)$ (hence $H_\theta(X) = H_0(X)$),

$$\begin{aligned} |I_\theta(X; Z) - I_0(X; Z)| &= |(H_\theta(Z) - H_0(Z)) - (H_\theta(X, Z) - H_0(X, Z))| \\ &\leq |H_\theta(Z) - H_0(Z)| + |H_\theta(X, Z) - H_0(X, Z)| \\ &\leq 2\left(\delta \log(|\mathcal{X}||\mathcal{Z}| - 1) + h_2(\delta)\right). \end{aligned}$$

Thus we may take

$$f(\varepsilon) := 2\left(\delta \log(|\mathcal{X}||\mathcal{Z}| - 1) + h_2(\delta)\right), \quad \delta := \sqrt{\frac{\varepsilon}{2}},$$

which satisfies $f(\varepsilon) \rightarrow 0$ as $\varepsilon \rightarrow 0$. \square

N. Decomposing Changes in Input Dependence

Definition N.1 (Entropy changes). Let X be prompts and let $Z \sim \pi_\theta(\cdot | X)$ under the current policy, with reference policy π_0 . Define the conditional-entropy and marginal-entropy changes

$$\Delta_{\text{in}} := H_\theta(Z | X) - H_0(Z | X), \quad \Delta_{\text{marg}} := H_\theta(Z) - H_0(Z).$$

The change in mutual information decomposes as $\Delta I = \Delta_{\text{marg}} - \Delta_{\text{in}}$. If an intervention (e.g., an entropy bonus) increases within-prompt variability $H(Z | X)$ more than it increases the marginal diversity $H(Z)$, input dependence necessarily decreases.

Theorem N.2. With Δ_{in} and Δ_{marg} defined above,

$$I_{\theta}(X; Z) - I_0(X; Z) = \Delta_{\text{marg}} - \Delta_{\text{in}}.$$

In particular, if $\Delta_{\text{in}} \geq \Delta_{\text{marg}} + \gamma$ for some $\gamma > 0$, then

$$I_{\theta}(X; Z) \leq I_0(X; Z) - \gamma,$$

and especially $I_{\theta}(X; Z) < I_0(X; Z)$ whenever $\Delta_{\text{in}} > \Delta_{\text{marg}}$.

Proof. Using $I(X; Z) = H(Z) - H(Z | X)$,

$$\begin{aligned} I_{\theta}(X; Z) - I_0(X; Z) &= (H_{\theta}(Z) - H_0(Z)) - (H_{\theta}(Z | X) - H_0(Z | X)) \\ &= \Delta_{\text{marg}} - \Delta_{\text{in}}. \end{aligned}$$

The sufficient-condition statements follow by rearranging the inequality. \square

An entropy bonus acts directly on the per-prompt dispersion and increases $H_{\theta}(Z | X)$, but it does not explicitly encourage cross-prompt separation that would increase the marginal entropy $H_{\theta}(Z)$ by a comparable amount. Hence it is plausible that Δ_{in} exceeds Δ_{marg} , in which case Theorem N.2 implies $I_{\theta}(X; Z)$ decreases.

Appendix L explains that when $\text{RV}(x)$ is small, the task update can be weak, so reward-agnostic regularizers can have larger relative influence on the total update.

O. GRPO Normalization Amplifies Noise at Low RV

GRPO-style normalization divides the advantage by $\sqrt{\text{RV}(x)}$, which induces a $\text{RV}(x)^{-1}$ noise amplification in the mean-squared error of the per-prompt gradient estimator.

For a fixed prompt x , define the normalized advantage

$$\tilde{A}(z; x) := \frac{A(z; x)}{\sqrt{\text{RV}(x)}}, \quad A(z; x) := R(z; x) - b(x), \quad b(x) := \mathbb{E}_{z \sim \pi_{\theta}(\cdot | x)}[R(z; x)].$$

Given K i.i.d. rollouts $z_1, \dots, z_K \sim \pi_{\theta}(\cdot | x)$, define

$$\hat{g}_{\text{GRPO}}(x) := \frac{1}{K} \sum_{k=1}^K \tilde{A}_k s_k, \quad g_{\text{GRPO}}(x) := \mathbb{E}[\tilde{A} s | X = x],$$

where $s_k = \nabla_{\theta} \log \pi_{\theta}(z_k | x)$.

Dividing the advantage by $\sqrt{\text{RV}(x)}$ causes the gradient estimator’s variance floor to scale as $\text{RV}(x)^{-1}$, so prompts with small reward variance suffer disproportionately noisy updates under GRPO-style normalization.

Proposition O.1 (GRPO variance floor). *Under Assumption I.1, the GRPO estimator satisfies*

$$\mathbb{E} \left[\|\hat{g}_{\text{GRPO}}(x) - g_{\text{GRPO}}(x)\|^2 | X = x \right] \geq \frac{1}{K} \cdot \frac{\sigma^2(x)}{\text{RV}(x)} \mathbb{E}[\|s\|^2 | X = x].$$

If $\sigma(x) = 0$, the lower bound is zero and thus vacuous.

This bound makes explicit that smaller $\text{RV}(x)$ yields a larger variance floor for the normalized estimator if all other factors are the same.

P. Core Author Contributions

Zihan Wang contributed across the full project lifecycle, including conceptualization, codebase and environment development, formal analysis, experiments, figures and plots, paper writing, and project correspondence. Chi Gui and Xing Jin contributed to the key ideas, software infrastructure, experiments, plots, and paper writing. Licheng Liu primarily contributed to the formal analysis and theory, experiments, and participated in paper writing. Qineng Wang contributed to the key ideas, software infrastructure, figures and plots, experiments, and paper writing.

Q. MI Proxy Family (extended table)

Table 6 below lists the full MI proxy family used in the main paper, including discrete and continuous variants.

Table 6. MI proxy family. All variants are derived from in-batch cross-scoring of reasoning traces against prompts, using matched (per-token log-prob under the true prompt) and marginal (per-token log-prob under the uniform prompt mixture) as base quantities. First-turn variants use only the first agent turn; trajectory variants sample across all turns.

Type	Proxy	Formula	Notes
Discrete	Retrieval-Acc	$\frac{1}{PG} \sum_{i,k} \mathbf{1}[\arg \max_j \mathbf{L}_{i,k,j} = i]$	Chance level $1/P$ under template collapse $k \in \{2, 4, 8\}$
	Recall@ k	$\frac{1}{PG} \sum_{i,k} \mathbf{1}[i \in \text{top-}k_j(\mathbf{L}_{i,k,j})]$	
Continuous (raw)	MI-Est	$\frac{1}{PG} \sum_{i,k} (\text{matched}_{i,k} - \text{marginal}_{i,k})$	Per-token; approaches 0 under collapse
	MI-Seq-Est	$\frac{1}{PG} \sum_{i,k} (\mathbf{L}_{i,k,i} - \log \frac{1}{P} \sum_j e^{\mathbf{L}_{i,k,j}})$	Per-sequence; no length normalization
Continuous (z-score)	MI-ZScore	$\frac{1}{PG} \sum_{i,k} \frac{\text{matched}_{i,k} - \text{marginal}_{i,k}}{\sigma_{\text{batch}} + \epsilon}$	Normalized by current-batch marginal std $\sigma_{\text{EMA}}^{(t)} = \alpha \sigma_{\text{EMA}}^{(t-1)} + (1-\alpha) \sigma_{\text{batch}}^{(t)}$
	MI-ZScore-EMA	$\frac{1}{PG} \sum_{i,k} \frac{\text{matched}_{i,k} - \text{marginal}_{i,k}}{\sigma_{\text{EMA}} + \epsilon}$	

Table 7. Sweep over prompt batch size P and group size G (trajectories per prompt) on Sokoban (Qwen2.5-3B). NF=no filtering; F=SNR-Aware Filtering ($\rho=0.9$). Total rollout budget is fixed at 128 trajectories across all configurations.

P (prompts) \times G (traj/prompt)	Task Perf. (%)			Step Time (s)			VRAM (GB)		
	NF	F	Δ	NF	F	$\Delta\%$	NF	F	Δ
128×1	23.6	–	–	89.8	–	–	201.80	–	–
64×2	18.8	27.3	+8.6	91.8	64.9	–29%	201.39	201.83	+0.44
32×4	24.2	27.4	+3.2	89.8	52.6	–41%	202.11	201.67	–0.44
8×16	15.6	23.6	+8.0	89.2	65.9	–26%	201.54	201.90	+0.36

R. Filtering Strategy Ablations

R.1. $P \times G$ rollout-budget sweep

Table 7 sweeps prompt batch size P and group size G on Sokoban (Qwen2.5-3B) at fixed rollout budget $K=128$. Configurations with $G \geq 4$ and SNR-Aware Filtering match or beat the 128×1 baseline at 26–41% lower per-step time.

Compute overhead of group sampling. SNR-Aware Filtering requires at least $G=2$ trajectories per prompt (group size) to estimate per-prompt RV. Since the total rollout budget is fixed at $K=128$ trajectories, varying the prompt batch size P and group size G is a repartitioning of that budget — all configurations incur identical rollout cost. Table 7 shows performance and wall-clock step time across configurations on Sokoban (Qwen2.5-3B). RV computation itself adds $<0.1\%$ of iteration time. With filtering ($\rho=0.9$), fewer groups enter gradient computation, reducing per-step time by 26–41%. Configurations with group size $G \geq 4$ and SNR-Aware Filtering match or outperform the 128×1 baseline, confirming that the gains come at no additional compute cost.

S. Stress-Testing the SNR Mechanism

The SNR framing makes a concrete causal claim: template collapse is a gradient-level consequence of low reward variance, not a side effect of aggressive regularization or model capacity. We stress-test this claim with four questions: (1) Does directly controlling RV level causally drive performance and MI? (2) Does injecting environmental noise predictably weaken MI? (3) Do gains come from signal quality rather than prompt-distribution bias? (4) When does the filtering condition hold in practice? A positive answer to all four makes the SNR account difficult to dismiss.

Quartile ablation provides direct causal evidence. To move beyond correlation between RV and performance, we run a controlled intervention. We sort all prompt groups by within-prompt RV, divide them into four quartiles (Q1 = highest, Q4 = lowest), and train four separate runs — each updating on one quartile only, all other settings fixed (Table 8). Task performance and MI degrade monotonically from Q1 to Q4. Combined with Theorem G.1 ($\|g_{\text{task}}\| \leq \sqrt{\text{RV}}$), this establishes the full causal chain: reward variance \rightarrow gradient quality \rightarrow input-dependent reasoning.

Controlled noise injection weakens MI. We run a direct intervention: varying environmental stochasticity and asking whether MI declines *predictably* in response. As environment and policy randomness increases, task return drops, conditional entropy rises, and $\hat{I}(X; Z)$ decreases monotonically (Figure 11). This is the expected consequence of the SNR chain. Additional noise inflates within-prompt return variance in a signal-free way, diluting the advantage estimates that task

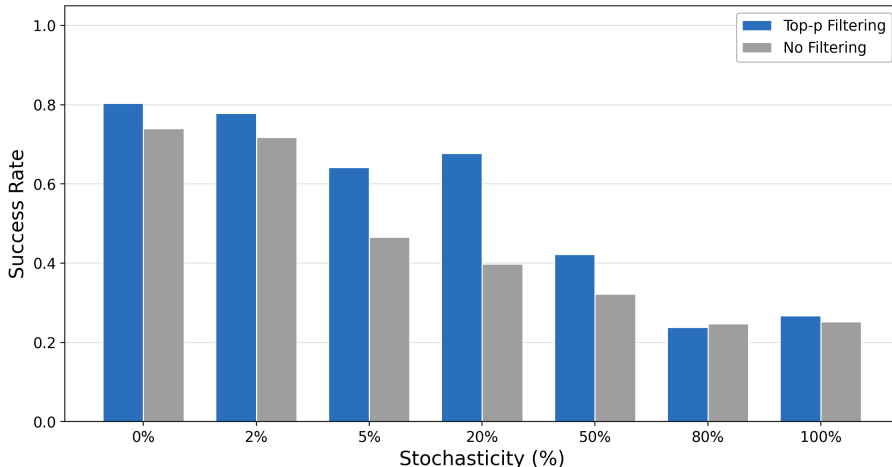


Figure 11. In FrozenLake, median success rates for both Top-p filtering (orange) and no filtering (gray) decrease as environment stochasticity increases from 0% to 100%. SNR-Aware Filtering maintains a clear advantage from 0% to 50% stochasticity, but the gap closes at 80%–100%, where high transition noise weakens reward variance as an informative signal proxy.

Table 8. Quartile ablation on Sokoban (Qwen2.5-3B, $P=8$, $G=16$, keeping 25% of prompts per step). Task performance and MI degrade monotonically from Q1 to Q4.

Quartile	RV Range	Task Perf (%)	MI Proxy	Entropy
Q1 (highest RV)	[4.4–5.6]	21.1	0.95	2.02
Q2	[1.5–4.2]	19.5	0.93	1.53
Q3	[0.0–0.2]	10.7	0.81	1.41
Q4 (lowest RV)	[0.0–0.1]	11.0	0.73	1.87

gradients depend on. Importantly, the filter’s advantage also attenuates at very high noise (80–100%), which is itself informative: when the environment is so stochastic that even high-effort prompts yield noisy rewards, RV loses its discriminative power. The mechanism predicts exactly this boundary condition.

Prompt-level filtering outperforms trajectory-level filtering. The gains from SNR-Aware filtering could come from selecting discriminative prompts, or from discarding hard/noisy trajectories. We disentangle these with a trajectory-level baseline: we keep all prompts but retain only the top-8 and bottom-8 trajectories per prompt by reward, preserving the prompt distribution while improving per-prompt SNR (Table 9). Trajectory-level filtering improves over no filtering. However, prompt-level SNR-Aware Filtering outperforms it by a wider margin. Within a naturally low-RV prompt, forcing within-prompt variance by sub-selecting trajectories amplifies noise. Selecting prompts that naturally produce discriminative signals is more effective.

Table 9. Trajectory-level vs. prompt-level filtering on Sokoban (Qwen2.5-3B). Prompt-level SNR-Aware Filtering provides the largest gains; trajectory-level filtering confirms that the benefit is not due to prompt-distribution bias.

Method	Prompts Used	Traj/Update	Task Perf (%)	MI Proxy
No filter	8/8	128	12.9	0.83
Prompt-level RV ($\rho=0.9$)	3.2/8	50.6	23.6	1.80
Trajectory-level	8/8	64	16.8	0.20

When does SNR-Aware Filtering help? Finally, we find SNR-Aware Filtering improves performance better when cross-prompt RV heterogeneity is large enough to separate signal-rich from noise-only prompts. We find the metric $\text{Std}(\text{RV})/\text{Mean}(\text{RV})$, computable from a single rollout batch, can effectively predict this (Table 10). When the ratio is high, the per-prompt RV distribution is bimodal and filtering cleanly separates signal from noise. When the ratio is near zero, all prompts carry similar RV and filtering discards data uniformly, like FrozenLake GRPO ($\Delta=-5.0\%$, ratio 0.33). This ratio is a cheap diagnostic which can be done before training.

How filtering adapts as training progresses? With the four predictions confirmed, we can now characterize how SNR-Aware Filtering behaves over the full training trajectory. Figure 12 tracks the effective kept ratio ρ_{eff} and zero-variance

Table 10. Per-setting RV statistics and filtering effectiveness. Std/Mean of RV predicts whether SNR-Aware Filtering helps: high ratio means bimodal RV and effective filtering; low ratio means uniform RV and random discarding.

Setting	Filter Δ	P	G	RV Mean	RV Std	RV Var	RV Min	RV Max	Std/Mean
Sokoban, 14B	+4.6%	8	8	2.24	2.88	8.32	0.10	6.00	1.29
Sokoban, 3B	+3.2%	32	4	2.49	2.89	8.35	0.05	6.52	1.16
FrozenLake, 3B (GRPO)	-5.0%	32	8	0.54	0.18	0.03	0.22	0.76	0.33

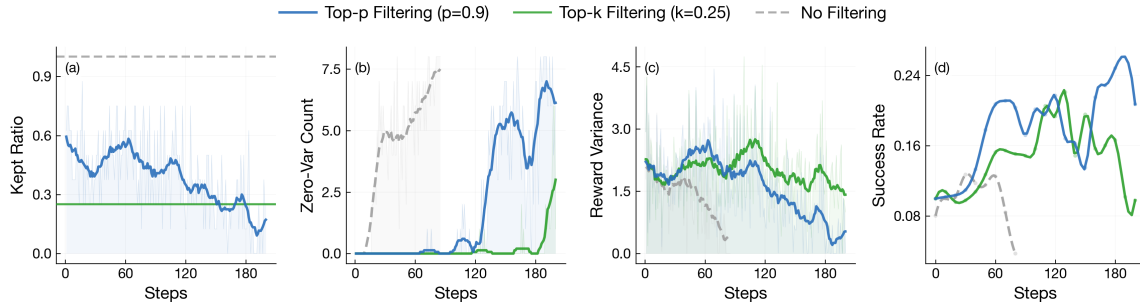


Figure 12. Effective kept ratio and zero-variance prompt count, showing adaptive selection pressure as variance collapses during training.

prompt count over training. Both move in the expected direction: as the policy improves and converges, more prompts yield near-identical rollout rewards (zero-variance count rises), and the filter responds by becoming more selective (kept ratio falls). This automatic tightening is precisely what a fixed strategy like Top- k with constant k cannot replicate. It would continue absorbing gradient budget from uninformative prompts even as signal quality deteriorates.

Reward collapse is visible at the distribution level. Figure 13 provides a complementary view of the same dynamics, tracking prompt-level reward distributions across early, mid, and late training in Sokoban. The shift is systematic: the hard portion shrinks as the policy improves, the mixed portion expands, and overall prompt-level variance collapses toward the late stages. This distribution-level signature mirrors the gradient-level story. Late training is not simply “easier” for the policy; it is a regime where reward variation has been compressed to the point that gradient updates carry progressively less task-discriminative information.

Format validity cannot substitute for content-sensitive diagnostics. One might hope that a coarser signal (whether the model’s output follows the required format) could serve as a collapse indicator without the overhead of MI estimation. Figure 14 shows this does not hold: format validity is largely decoupled from collapse, with runs maintaining near-perfect validity while exhibiting low MI. Structural correctness and semantic input-dependence are separate dimensions. This reinforces the need for content-sensitive diagnostics, and explains why the MI proxy provides signal that format-based checks miss.

RV is largely orthogonal to entropy and response length, which explains why entropy-based stabilizers cannot prevent template collapse. Reward variance correlates weakly with conditional entropy (Spearman -0.14) and response length (0.12), while correlating strongly with task reward (0.63). RV therefore targets a distinct axis of update quality rather than surface statistics, making it a complementary control knob to KL and entropy regularization. Figure 12 further shows that the effective kept ratio adapts over training: as more prompts drift toward near-zero RV, the filter automatically concentrates gradient updates on the shrinking pool of still-informative prompts.

What is the relationship between SNR-Aware Filtering and KL/entropy tuning stabilization? When training RL agents, practitioners typically tune KL penalty and entropy regularization coefficients to maintain training stability and prevent mode collapse. However, these interventions primarily control within-input diversity ($H(Z | X)$) and cannot directly address the signal-to-noise imbalance that drives template collapse. Even with carefully tuned regularization, if most prompts have low reward variance, the task gradient remains weak and regularization forces still dominate the update direction.

SNR-Aware Filtering is complementary: it selects high-signal prompts at each iteration, directly boosting the fraction of task-discriminative gradient in each update. This acts as a signal-enhancement mechanism rather than a noise-control mechanism. We provide a detailed empirical comparison of KL tuning, entropy tuning, and SNR-Aware Filtering in Section 5.1, showing that the three interventions move training dynamics along different axes (Figure 8).

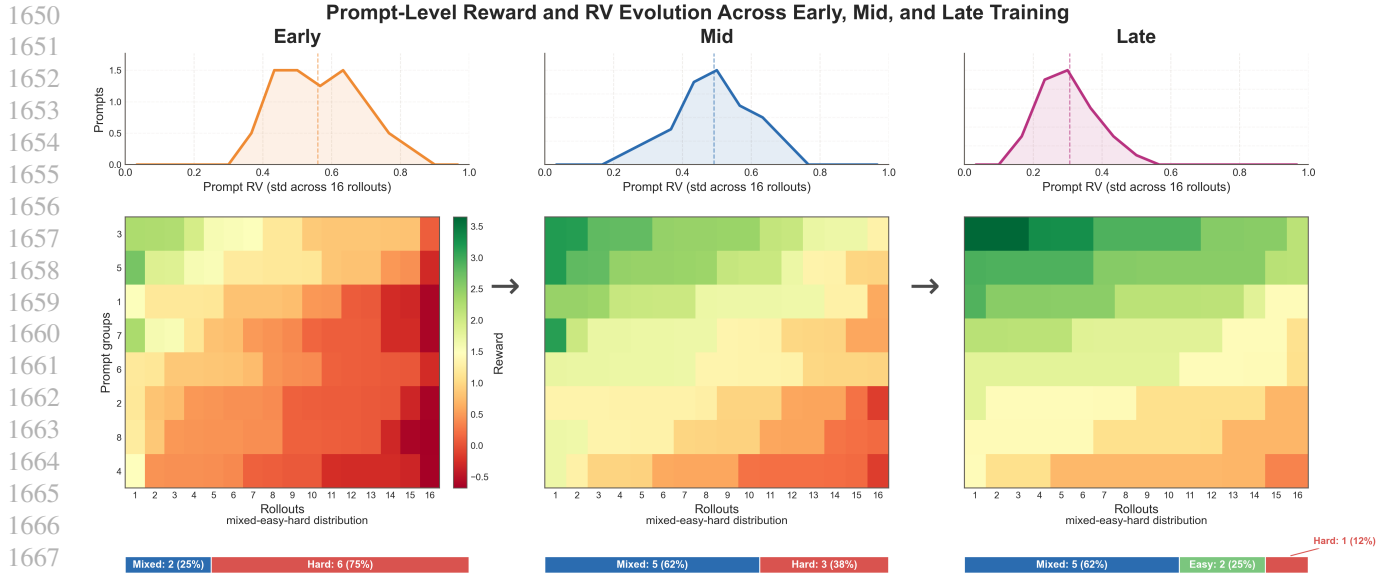


Figure 13. Prompt-level reward distribution across training phases, showing RV collapse as prompts shift toward uniform reward structures.

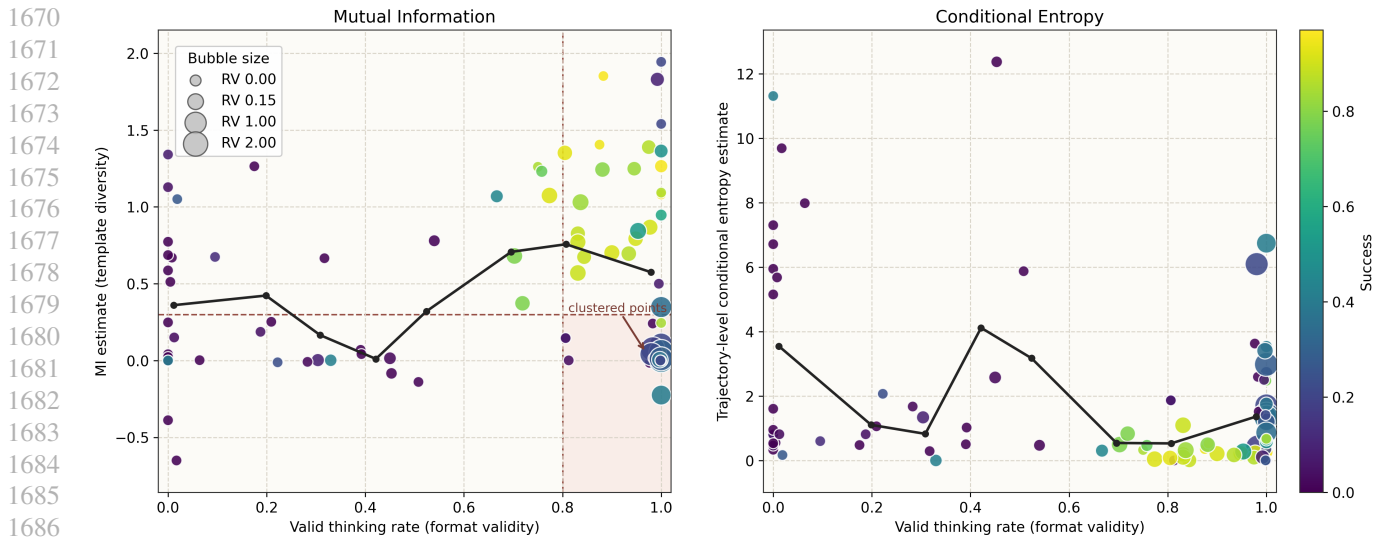


Figure 14. Format validity versus MI and entropy diagnostics, showing that high validity does not guarantee high input dependence.

T. Extended Limitations

The SNR decomposition assumes task-signal and regularization noise separate cleanly, though they may couple through gradient accumulation in practice. All experiments are single-agent; how template collapse propagates in multi-agent RL remains open. A capable model could game the filtering criterion by artificially inflating reward variance, a risk worth monitoring over long training horizons. The method requires reward variance to be a reliable signal proxy, which degrades in sparse or noisy reward environments. Aggressive filtering may narrow exploration coverage; the kept mass requires per-task tuning.