ENTROPY-PRESERVING REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Policy gradient algorithms have been a driver of much recent advancement in language model reasoning. One of their most appealing properties is the ability to learn from exploration on their own trajectories, a process crucial for discovering diverse approaches and fostering creative solutions. As we show in this paper, most policy gradient algorithms naturally reduce the entropy—and thus the diversity of explored trajectories—as part of training, yielding a policy increasingly limited in its ability to explore. However, not all algorithms exhibit this collapse in entropy equally. In this paper, we formally analyze the contributions of leading policy gradient objectives on entropy, show which mechanisms they employ to implicitly limit entropy collapse, and propose a new regularization method, REPO, that stabilizes entropy over training through the use of an adaptive controller. Models trained with REPO preserve entropy throughout training, yielding final policies that are, on average, more performant. By preserving entropy in the final policy, REPO-trained models can even be re-trained on evolving data distributions in new environments, unlike their non-entropy-preserving counterparts.

1 Introduction

Online policy gradient reinforcement learning (RL) has become the standard for boosting the reasoning abilities of language models (Jaech et al., 2024; Comanici et al., 2025; Guo et al., 2025). This approach involves sampling trajectories from the current policy within a given environment and reward function, then using these to estimate a gradient that maximizes expected reward. Effective RL optimization requires balancing exploration and exploitation (Thrun, 1992; Sutton et al., 1998), where a robust learner should generate diverse trajectories to cover the spectrum of potential solutions. Maximum entropy reinforcement learning offers a framework for achieving this balance (Ziebart et al., 2008; Haarnoja et al., 2017; 2018; Eysenbach & Levine, 2022). While trivially the optimal solution to a finite Markov decision process (MDP) is a deterministic stationary policy, optimization over the intermediate landscape requires a balance of exploration and exploitation.

A common issue observed in online algorithms like GRPO (Shao et al., 2024) is entropy collapse. This phenomenon occurs when training excessively narrows the distribution around already high-probability solutions from the base model, neglecting other correct but less probable options. This often yields premature convergence to a local optimum, enhancing pass@1 relative to base model at the expense of pass@k (Shao et al., 2024; Dang et al., 2025; Yue et al., 2025). This challenge has spurred innovations in policy gradient algorithm design, e.g. directly optimizing for pass@k performance (Chen et al., 2025b). Concurrently, research has highlighted GRPO's training instability and the complex interplay between off-policy drift, importance weight clipping, and entropy, inspiring modifications such as DAPO (Yu et al., 2025) and GSPO (Zheng et al., 2025).

In this work, we analyze entropy preservation as a unifying lens for understanding the successes of recent algorithms and to propose a novel family of policy gradient objectives. An important observation from our work is that, while a correlation exists between final entropy and performance, a more informative measure is the entropy trajectory throughout the optimization process. As the saying goes, "it's not the destination, it's the journey." Figure 1 tracks this effect. A trajectory characterized by lower entropy throughout yields lower performance. Conversely, if entropy trajectories are similar for most of the optimization but differ only in the final steps, performance is largely unaffected.

Given this observation, we turn to study the entropy behavior of various leading RL algorithms. We begin by theoretically analyzing how the REINFORCE policy gradient objective modulates entropy

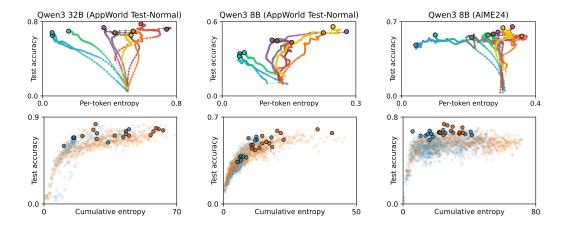


Figure 1: Top: Evolution of the average per-token entropy and test accuracy during training for several baseline (GRPO, LOOP, DAPO, GSPO) and their entropy regularized versions (REPO). Each curve shows the average trajectory over several training runs with different seeds. Bottom: Cumulative entropy experienced during training up to a given checkpoint is positively correlated with the test accuracy. Each point is a checkpoint of a single training run (best-performing checkpoint per run highlighted). Algorithms that collapse the entropy early (see Qwen-3-8B on AppWorld; middle column) perform significantly worse than algorithms that maintain a steady entropy during training. See App. C for a detailed study and breakdown of this phenomena across algorithms.

dynamics, explaining how entropy can decrease during training. This effect is amplified when talking multiple policy update steps, typical of GRPO and other PPO-like algorithms. We also show how importance weight clipping and its modifications as seen in DAPO and GSPO, can mitigate this pressure. Finally, we show that regularizing the entropy during training allows a broad family of policy gradient algorithms to train more performant policies. Specifically, we introduce Regulated Entropy Policy Optimization (REPO), an approach to policy gradient optimization that adaptively reweighs advantages and log-probabilities online to preserve entropy. REPO uses an adaptive controller, tracking entropy dynamics live, and adjusting regularization strength accordingly. Training with REPO achieves state-of-the-art results on AppWorld and strong performance on AIME 2024 and AIME 2025. Furthermore, we demonstrate that policies trained with REPO retain their trainability, allowing for iterative learning on new tasks in novel environments, a capability often lost in policies trained without explicit entropy preservation.

2 PRELIMINARIES

Language modeling. Let $x \in \mathcal{X}$ denote the tokens in a vocabulary and $x \in \mathcal{X}^*$ the strings expressible via concatenation of those tokens. A language model (LM) π_{θ} parameterized by θ defines a probability distribution over strings that factors autoregressively such that $\pi_{\theta}(x) = \pi_{\theta}(\Box \mid x) \prod_{i=1}^{|x|} \pi_{\theta}(x_i \mid x_{< i})$, where \Box denotes an end of sequence (EOS) marker. Note that for notational convenience we will use π_{θ} to express probabilities on both tokens and strings.

Language modeling as a Markov decision process. Let the policy π_{θ} sample actions $a \in \mathcal{A} = \mathcal{X} \cup \{\Box\}$ (any token or EOS) given a **state** $s \in \mathcal{X}^*$ (a string context). Let state transitions append generated actions to the state. Let τ denote a **trajectory**, a sequence of states and actions generated by the policy and environment. Let $\tau \sim \pi_{\theta}$ denote the trajectory distribution. We consider tasks with **terminal rewards** $R(c, \tau)$. Given some task context c sampled from some dataset c, the MDP objective is to maximize $\mathcal{J}_{\text{MDP}} \stackrel{\text{def}}{=} \mathbb{E}_{c \sim \mathcal{D}, \tau \sim \pi_{\theta}(\cdot | c)}[R(c, \tau)]$.

¹State transitions deterministically append the generated action to the context, terminating generation at EOS or upon some other environment condition. In some domains, e.g., those involving tool calls, state transitions may also append additional tokens to the state that were generated by some unobservable process such as executing a code interpreter.

Policy gradient reinforcement learning directly computes a gradient through the REINFORCE algorithm (Williams, 1992), which is amenable to Monte Carlo estimation:

$$\nabla_{\theta} \mathcal{J}_{\text{MDP}} = \mathbb{E}_{\boldsymbol{c} \sim \mathcal{D}, \boldsymbol{\tau} \sim \pi_{\theta}(\cdot | \boldsymbol{c})} \left[A(\boldsymbol{c}, \boldsymbol{\tau}) \cdot \nabla_{\theta} \log \pi_{\theta}(\boldsymbol{\tau} | \boldsymbol{c}) \right],$$

where $A(c, \tau) = R(c, \tau) - b$ is an advantage function shifting the return $R(c, \tau)$ by a baseline b.

REINFORCE leave-one-out (RLOO) (Kool et al., 2019; Ahmadian et al., 2024; Kazemnejad et al., 2024; Chen et al., 2025a) is one of the most popular estimates of advantage for language modeling. It generates K independent samples on-policy $\tau_1, ..., \tau_k \sim \pi_\theta(\cdot \mid \mathbf{c})$ for each task \mathbf{c} . The reward for each trajectory may then be baselined against the remaining K-1 independent samples, yielding an unbiased, low variance advantage estimator:

$$\widehat{A}_{\text{RLOO}}(\boldsymbol{c}, \boldsymbol{\tau}_i) \stackrel{\text{def}}{=} R(\boldsymbol{c}, \boldsymbol{\tau}_i) - \frac{1}{K-1} \sum_{j=1}^K R(\boldsymbol{c}, \boldsymbol{\tau}_j) \mathbb{1}_{[i \neq j]} = \frac{K}{K-1} \left(R(\boldsymbol{c}, \boldsymbol{\tau}_i) - \frac{1}{K} \sum_{j=1}^K R(\boldsymbol{c}, \boldsymbol{\tau}_j) \right).$$

Policy gradient algorithms are on-policy by nature: They rely on a new set of trajectories in each context $\tau \sim \pi_{\theta}(\cdot \mid c)$ after each gradient update of the policy π_{θ} .

Proximal policy optimization (PPO) allows the updated policy to deviate slightly from a sampling policy (Schulman et al., 2017). It uses an importance weight to correct the magnitudes of parameter updates such that the expected policy gradient remains unbiased. These importance weights are typically clipped to avoid divergence from a local trust region (Schulman et al., 2015).

$$\mathcal{J}_{\text{PPO}} \stackrel{\text{def}}{=} \mathbb{E}_{\boldsymbol{c} \sim \mathcal{D}, \boldsymbol{\tau} \sim \pi_{\theta}(\cdot \mid \boldsymbol{c})} \left[\frac{1}{|\boldsymbol{\tau}|} \sum_{a_{t} \in \boldsymbol{\tau}} \min \left(A(\boldsymbol{c}, \boldsymbol{\tau}) w_{t}, A(\boldsymbol{c}, \boldsymbol{\tau}) w_{t} |_{1-\epsilon}^{1+\epsilon} \right) \right] \quad w_{t} \stackrel{\text{def}}{=} \frac{\pi_{\theta}^{\text{new}}(a_{t} \mid \boldsymbol{c}, \boldsymbol{a}_{< t})}{\pi_{\theta}^{\text{old}}(a_{t} \mid \boldsymbol{c}, \boldsymbol{a}_{< t})},$$

where $w_t|_{1-\epsilon}^{1+\epsilon}$ clips the importance ratio from below $1-\epsilon$ and above $1+\epsilon$. In our theoretical analysis, we will examine PPO with and without clipping. The version studied will be clear from the context.

LOOP (Chen et al., 2025a) and **GRPO** (Shao et al., 2024) combine the above PPO objective with RLOO leave-one-out advantage estimates. GRPO rescales advantages by the standard deviation of the sample returns, which introduces a small bias (Liu et al., 2025b),

$$\widehat{A}_{\mathrm{GRPO}}(\boldsymbol{c}, \boldsymbol{\tau}_i) \stackrel{\mathrm{def}}{=} \frac{R(\boldsymbol{c}, \boldsymbol{\tau}_i) - \mathrm{mean}(R(\boldsymbol{c}, \boldsymbol{\tau}_1)), \ldots, R(\boldsymbol{c}, \boldsymbol{\tau}_k))}{\mathrm{std}(R(\boldsymbol{c}, \boldsymbol{\tau}_1)), \ldots, R(\boldsymbol{c}, \boldsymbol{\tau}_k))}$$

while LOOP uses \widehat{A}_{RLOO} directly.

Group Sequence Policy Optimization (GSPO) Zheng et al. (2025) uses a trajectory-level trust region defined by the geometric average of a sequence's probability ratios

$$\mathcal{J}_{\text{GSPO}} \stackrel{\text{def}}{=} \mathbb{E}_{\boldsymbol{c} \sim \mathcal{D}, \boldsymbol{\tau} \sim \pi_{\theta}(\cdot | \boldsymbol{c})} \left[\min \left(A(\boldsymbol{c}, \boldsymbol{\tau}) w^{\text{GSPO}}, A(\boldsymbol{c}, \boldsymbol{\tau}) w^{\text{GSPO}} \big|_{1 - \epsilon}^{1 + \epsilon} \right) \right] \quad w^{\text{GSPO}} \stackrel{\text{def}}{=} \left(\frac{\pi_{\theta}^{\text{new}}(\boldsymbol{\tau} \mid \boldsymbol{c})}{\pi_{\theta}^{\text{old}}(\boldsymbol{\tau} \mid \boldsymbol{c})} \right)^{\frac{1}{|\boldsymbol{\tau}|}}.$$

GSPO yields an equivalent gradient estimator to GRPO, LOOP, and RLOO on-policy, but clips tokens and trajectories differently as the updated policy $\pi_{\theta}^{\text{new}}$ drifts from the sampling policy $\pi_{\theta}^{\text{old}}$.

Policy entropy. The inherent uncertainty that a policy places over its generations may be expressed from an information theoretic standpoint as **entropy** – expected surprise: $\mathcal{H}_{\pi_{\theta}}(\mathcal{D}) = -\mathbb{E}_{\boldsymbol{c} \sim \mathcal{D}}\left[\mathbb{E}_{\boldsymbol{\tau} \sim \pi_{\theta}(\cdot | \boldsymbol{c})}\left[\log \pi_{\theta}(\boldsymbol{\tau} \mid \boldsymbol{c})\right]\right]$. In addition to global entropy, we may consider the entropy over actions at any given state $\boldsymbol{s} = (\boldsymbol{c}, \boldsymbol{a}_{< t})$ as $\mathcal{H}_{\pi_{\theta}}(\cdot \mid \boldsymbol{s}) = -\mathbb{E}_{a_{t} \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})}\left[\log \pi_{\theta}(a_{t} \mid \boldsymbol{s})\right]$.

In this paper, we show how state-wise entropy changes as variants of policy gradient optimize their objectives. We show which variants are naturally entropy preserving, and which variants lead to a rapid collapse. Finally, we show that a simple class of transformations applied to the advantages lead to a very simple and effective regularization of entropy.

3 THE ENTROPY DYNAMICS OF POLICY GRADIENT

The entropy dynamics of policy gradient RL boils down to the relationship between two values: (1) action log-probabilities, and (2) the advantages yielded by those actions. Intuitively, assigning

a positive advantage to some action increases its probability. For high probability actions, this effect sharpens the distribution, and entropy decreases. For low probability actions, this flattens the distribution, increasing entropy. The opposite pattern holds for negative advantages. This effect is quite natural. After all, sharpening an uncertain policy around correct actions directly maximizes the expected return. However, as we will see, not all RL algorithms sharpen the distribution equally.

Formally, consider the policy gradient update with on-policy actions in state s. Under a first-order Taylor approximation to the training dynamics, the expected change in entropy is as follows.

Theorem 1. Given a policy gradient update $\hat{\theta} := \theta + \alpha \cdot \nabla_{\theta} \mathcal{J}_{MDP}(s)$, the expected change in entropy is approximately:

$$\Delta \mathcal{H}_{\pi_{\theta}}(\cdot \mid \boldsymbol{s}) \approx -\alpha \cdot \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s}), a' \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[A(\boldsymbol{s}, a) \cdot \left(\log \pi_{\theta}(a' \mid \boldsymbol{s}) - \overline{\log \pi_{\theta}(\cdot \mid \boldsymbol{s})} \right) \cdot u(\boldsymbol{s}, a)^{\top} u(\boldsymbol{s}, a') \right],$$

where $\overline{\log \pi_{\theta}(\cdot \mid \mathbf{s})} = \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s})}[\log \pi_{\theta}(a \mid \mathbf{s})]$ and $u(\mathbf{s}, a) \stackrel{\text{def}}{=} \nabla_{\theta} \log \pi_{\theta}(a \mid \mathbf{s})$ is the score function of π_{θ} evaluated at state \mathbf{s} and action a.

[Proof in App. A.2]. The entropy change is driven by a multiplicative relationship between action log-probabilities and the advantages yielded by those actions. In an exact derivation, these are weighted by the score vector outer product. With additional independence assumptions or a tabular softmax policy parameterization, this expression can be further simplified, resulting in a weighting by the action probabilities, which may explain various behaviors of popular algorithms. This yields the following corollary:

Corollary 1. Assuming $u(s, a)^{\top}u(s, a') = 0$ for all $a \neq a'$, the change in entropy is approximately:

$$\Delta \mathcal{H}_{\pi_{\theta}}(\cdot \mid \boldsymbol{s}) \propto -\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[A(\boldsymbol{s}, a) \cdot \left(\log \pi_{\theta}(a \mid \boldsymbol{s}) - \overline{\log \pi_{\theta}(\cdot \mid \boldsymbol{s})} \right) \cdot \pi_{\theta}(a \mid \boldsymbol{s}) \right]$$

[Proof in App. A.3]. This latter form encodes the dominant behavior of entropy dynamics in a manner that is inherent to policy gradient. Using this form, we explain the observed behaviors of various RL algorithms. A similar derivation can be shown for tabular softmax policies (Cui et al., 2025, see Corollary 2 in App. A.4). Corollaries 1 and 2 tell us that the change in entropy is governed by a correlation between advantages and log-probabilities, weighted by action probability.

Entropy dynamics of PPO. The biggest feature of PPO is its ability to train on slightly off-policy trajectories, given that the updated policy does not deviate from a trust region around the current policy. This allows PPO to take multiple policy-improvement steps for a single set of trajectories. The effect of these repeated updates are much larger policy updates between consecutive PPO steps, which empirically amplifies entropy collapse. This being said, the clipping on PPO, when appropriately orchestrated, can protect against entropy collapse as well. Clipping ensures that no policy gradient update is performed if the policy drifts outside a trust region $(1 - \epsilon_{\text{low}}) \cdot \pi_{\theta}^{\text{old}}(a \mid s) \leq \pi_{\theta}^{\text{new}}(a \mid s) \leq (1 + \epsilon_{\text{high}}) \cdot \pi_{\theta}^{\text{old}}(a \mid s)$. This bounds the change in entropy:

Theorem 2. Proximal Policy Optimization (PPO) bounds the entropy $\mathcal{H}_{\pi_{\theta}^{\text{new}}}(\cdot \mid s)$ of the updated policy by the original policy entropy $\mathcal{H}_{\pi_{\theta}^{\text{old}}}(\cdot \mid s)$ such that:

$$(1 - \epsilon_{low}) \cdot \mathcal{H}_{\pi_{\theta}^{old}}(\cdot \mid s) \leq \mathcal{H}_{\pi_{\theta}^{new}}(\cdot \mid s) \leq (1 + \epsilon_{high}) \cdot \mathcal{H}_{\pi_{\theta}^{old}}(\cdot \mid s)$$

[Proof in App. A.5]. The clipping thresholds directly limit the maximum induced change in entropy per token. Intuitively, the change in entropy per token is stochastic: some actions have a large correlation between advantage and log probability; others do not, or even have an anti-correlation. For a symmetric clipping regime, this results in an entropy change that largely follows the statistical trends outlined above, but at a lower magnitude.

Entropy dynamics of DAPO. Now consider DAPO (Yu et al., 2025), with an asymmetric clipping regime $\epsilon_{\text{low}} < \epsilon_{\text{high}}$. This allows for larger entropy increases, while limiting the entropy decrease. Due to the stochastic nature of the entropy changes, this directly contributes to an overall increase in per-token entropy over sufficient samples. Threshold values $\epsilon_{\text{low}} = 0.2$ and $\epsilon_{\text{high}} = 0.28$ proposed in Yu et al. (2025) stabilize the entropy throughout training, as we show experimentally.

Entropy dynamics of GSPO. GSPO defines a trust region $1 - \epsilon_{\rm low}^{\rm GSPO} \leq w^{\rm GSPO} \leq 1 + \epsilon_{\rm high}^{\rm GSPO}$, or equivalently $\left(1 - \epsilon_{\rm low}^{\rm GSPO}\right)^{|\tau|} \leq \frac{\pi_{\theta}^{\rm new}(\tau|c)}{\pi_{\theta}^{\rm old}(\tau|c)} \leq \left(1 + \epsilon_{\rm high}^{\rm GSPO}\right)^{|\tau|}$. This induces an equivalent bound to Thm. 2; however, the bound now depends on the trajectory length $|\tau|$. Longer trajectories may induce a larger change in entropy, shorter trajectories induce a smaller change in entropy. With parameter values suggested in Zheng et al. (2025), $\epsilon_{\rm low}^{\rm GSPO} = 3 \times 10^{-4}$ and $\epsilon_{\rm high}^{\rm GSPO} = 4 \times 10^{-4}$, the entropy bound is tighter for trajectories $|\tau| < \frac{\ln(1\pm\epsilon)}{\ln(1\pm\epsilon^{\rm GSPO})} \approx 600$ tokens compared to DAPO. Like DAPO, the clipping range is asymmetric $\epsilon_{\rm low}^{\rm GSPO} < \epsilon_{\rm high}^{\rm GSPO}$ leading to a stochastic increase in entropy.

Regulated entropy policy optimization (REPO) changes the advantage function to include a scaled policy log-likelihood term $A_{\text{REPO}}(s,a) = A(c,a) - \beta_s(\log \pi_\theta(a_t|s) - \overline{\log \pi_\theta(\cdot \mid s)})$ for each $s = (c, a_{< t})$. This updated advantage is no longer constant throughout the trajectory like in RLOO and variants, but differs for individual tokens $a_t \in \tau$. Following Thm. 1 by Prop. 3 the induced change in entropy with A_{REPO} is:

$$\Delta \mathcal{H}_{\pi_{\theta}}^{\text{REPO}}(\cdot \mid \boldsymbol{s}) \approx \Delta \mathcal{H}_{\pi_{\theta}}(\cdot \mid \boldsymbol{s}) + \beta_{\boldsymbol{c}} \cdot \underbrace{\alpha \cdot \left\| \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[(\log \pi_{\theta}(\boldsymbol{s}, a) - \overline{\log \pi_{\theta}(\cdot \mid \boldsymbol{s})}) \cdot u(\boldsymbol{s}, a) \right] \right\|^{2}}_{>0}.$$

This provides us with a direct mechanism to control the entropy. A positive $\beta_c > 0$ increases the entropy over actions in a state relative to the default dynamic. A value $\beta_c = 0$ allows the natural entropy decrease to proceed. A negative $\beta_c < 0$ collapses the entropy. Note, this holds for any parametrization of the policy and does not rely on approximations.

How should we choose β_c to preserve entropy? One natural choice is to counter-act the entropy collapse on a per-token level and set $\beta_c^{\text{REPO-D}} \propto -\Delta \mathcal{H}_{\pi_\theta}(\cdot \mid s)$ as approximated in Corollary 1, thus neutralizing $\Delta \mathcal{H}_{\pi_\theta}$ and allowing $\Delta \mathcal{H}_{\pi_\theta}^{\text{REPO}}$ to approach 0 with the right scale of the regularizer. We call this variant REPO-D.

While the above heuristic provides us with an overall mechanism to control entropy, the exact scale of the regularizer depends on many aspects of the policy gradient optimization: the learning rate, the structure of the gradient, second order effects, etc. We learn the magnitude ζ of the regularizer using a simple control heuristic similar to the adaptive $D_{\rm KL}$ penalty presented in Schulman et al. (2017). Let $\beta_c^{\rm REPO-D} = \zeta \Delta \mathcal{H}_{\pi_\theta}(\cdot \mid s)$. The heuristic proceeds as follows: (1) Estimate $\mathcal{H}_{\pi_\theta}^{\rm init}$, the policy entropy over the experience collected in this first iteration. (2) Initialize $\zeta = 10^{-3}$. (3) On each iteration, estimate \mathcal{H}_{π_θ} , the current policy entropy, and compare it to $\mathcal{H}_{\pi_\theta}^{\rm init}$. If $\mathcal{H}_{\pi_\theta} < \mathcal{H}_{\pi_\theta}^{\rm init}$, update $\zeta \leftarrow \zeta \times 2$. If $\mathcal{H}_{\pi_\theta} > \mathcal{H}_{\pi_\theta}^{\rm init}$, update $\zeta \leftarrow \zeta \div 2$. (4) Clip ζ to the window $\zeta_{\rm min} \leq \zeta \leq \zeta_{\rm max}$.

Supporting rare correct actions (REPO-R). Looking back through our learnings thus far, it appears that the most important bang-for-buck in preserving entropy is through raising low probability correct actions. This intuitively corresponds to reinforcing rare but correct solutions under our policy optimization, which is a behavior that we hope to encourage. We can thus apply an entropy regularizer on positive advantage actions only: $\beta_{a,c}^{\rm REPO-R} = -\zeta \max(A(c,a),0)$. The effect of this is simple: Wrong (negative-advantage) actions are unaffected and penalized by the negative advantage. For correct actions (positive-advantage) the entropy regularizer reduces advantages for high-probability outcomes, but amplifies low-probability samples (towards a higher entropy state). This does however introduce a small bias to the gradient estimate, as it treats positive and negative advantage samples differently. Different scales of regularizers $\beta_c^{\rm REPO-D}$ and $\beta_{a,c}^{\rm REPO-R}$ demand different clipping ranges $[\zeta_{\rm min}, \zeta_{\rm max}]$: $[10^{-3}, 10^{1}]$ for REPO-D and $[10^{-5}, 10^{-1}]$ for REPO-R.

4 EXPERIMENTS

With the theory established, we evaluate whether training with REPO yields improvements to strong models on challenging environments when compared to state-of-the-art learning algorithms. We choose Qwen-3-8B and Qwen-3-32B as our starting policies (Yang et al., 2025).

Environments. Interactive tool-use agent. Training scenarios are drawn from the train split (90 problems) of the AppWorld benchmark (Trivedi et al., 2024). The AppWorld Test Normal (TN,

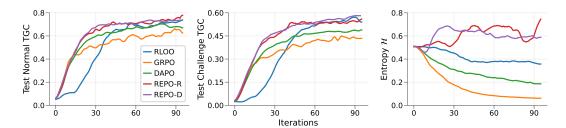


Figure 2: Qwen-3-32B AppWorld test performance and token entropy across iterations of training. Each curve shows a mean across various random seeds.

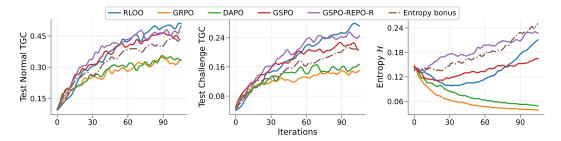


Figure 3: Qwen-3-8B AppWorld test performance and token entropy across iterations of training. Each curve shows a mean across various random seeds.

168 tasks) and Test Challenge (*TC*, 417 tasks) splits are used for evaluation. Terminal reward is calculated via task-provided unit-tests that check the final state of the environment against ground truth (additional details in App. B.1). *Competition-level mathematics*. Training scenarios are drawn from a non-overlapping quality-filtered subset of the AMC/AIME section of NuminaMath-1.5 (563 problems; Li et al., 2024). AIME 2024 (30 problems) and AIME 2025 (30 problems) are used as evaluation datasets. Terminal reward indicates whether the generated answer matches the reference.

Training details. We run training for 100 iterations on AppWorld, and 200 on AIME, where each iteration consists of an experience collection stage followed by policy optimization on the collected data. For each iteration, 64 (AppWorld) or 128 (AIME) scenarios are sampled without replacement and K=6 rollouts are generated per scenario. Model outputs are generated with the thinking template enabled using temp = 1.0 and a maximum generation length of 16384 (AppWorld) or 4096 (AIME) tokens.² All scenarios yielding identical in-group returns ($\hat{A}_{\rm RLOO}=0$) are filtered out to increase throughput. Additional training details in App. B.2.

Algorithms. For each algorithm, we highlight its distinguishing features with otherwise minimal deviations from the base policy gradient to aid reproducibility (thus, some details and hyperparameter choices may differ slightly from original sources).

RLOO: REINFORCE with the $\widehat{A}_{\text{RLOO}}$ advantage estimator. The collected experience is trained strictly on-policy for 1 epoch using a large minibatch that comprises all collected experience. *GRPO*: Off-policy extension of RLOO that uses normalized Leave None Out (LNO) estimator $\widehat{A}_{\text{GRPO}}$ and symmetric PPO clipping with $\epsilon=0.2$. With GRPO and all algorithms below we train on the collected experience for 2 epochs with the minibatch size of 128 (AppWorld) or 256 (AIME) trajectories. *LOOP*: A variant of GRPO with non-normalized estimator $\widehat{A}_{\text{RLOO}}$ (RL SOTA on AppWorld at the time of writing). DAPO: a variant of LOOP with asymmetric clipping ($\epsilon_{\text{low}}=0.2$, $\epsilon_{\text{high}}=0.28$). *GSPO*: An adjustment of LOOP using w^{GSPO} importance weighting and trajectory-based clipping with $\epsilon_{\text{low}}^{\text{GSPO}}=3\times10^{-4}$ and $\epsilon_{\text{high}}^{\text{GSPO}}=4\times10^{-4}$. *REPO* algorithms incorporate entropy control over either DAPO (e.g. in REPO-R) or GSPO (in GSPO-REPO-R) as their base algorithm.

 $^{^2}$ This restricted context window was chosen for experiment iteration speed. Note that Qwen-3 models can be trained to higher accuracy on AIME with increased token budget.

Algorithm	Test Normal	Best TN	Test Challenge	Best TC	\mathcal{H}	$\Delta \mathcal{H}$
RLOO	0.71 ± 0.02	0.73	0.52 ± 0.02	0.54	0.32	-0.19
GRPO	0.61 ± 0.01	0.62	0.42 ± 0.03	0.45	0.06	-0.45
LOOP	0.64 ± 0.04	0.67	0.40 ± 0.01	0.41	0.05	-0.46
DAPO	0.70 ± 0.03	0.73	0.48 ± 0.01	0.49	0.18	-0.33
GSPO	0.62 ± 0.06	0.68	0.51 ± 0.05	0.56	0.42	-0.09
REPO-R	0.73 ± 0.03	0.75	0.55 ± 0.05	0.62	0.48	-0.03
REPO-D	0.71 ± 0.03	0.74	0.54 ± 0.02	0.55	0.63	+0.12
GSPO-REPO-R	0.71 ± 0.00	0.71	0.54 ± 0.03	0.57	0.57	+0.06
GSPO-REPO-D	0.64 ± 0.02	0.66	0.41 ± 0.02	0.43	0.52	+0.01

Table 1: Task goal completion scores for AppWorld Qwen-3-32B by training algorithm.

Algorithm	Test Normal	Best TN	Test Challenge	Best TC	H	$\Delta \mathcal{H}$
RLOO	0.48 ± 0.07	0.54	0.25 ± 0.02	0.28	0.14	-0.00
GRPO	0.33 ± 0.01	0.34	0.14 ± 0.01	0.15	0.04	-0.10
LOOP	0.27 ± 0.01	0.28	0.14 ± 0.00	0.14	0.04	-0.10
DAPO	0.34 ± 0.02	0.36	0.17 ± 0.01	0.17	0.05	-0.09
GSPO	0.45 ± 0.01	0.46	0.21 ± 0.01	0.21	0.15	+0.01
${\cal H}$ Bonus	0.42 ± 0.04	0.45	0.20 ± 0.01	0.22	0.22	+0.08
REPO-R	0.43 ± 0.01	0.45	0.21 ± 0.00	0.21	0.15	+0.01
REPO-D	0.43 ± 0.00	0.43	0.19 ± 0.01	0.20	0.14	-0.00
GSPO-REPO-R	0.48 ± 0.04	0.51	0.27 ± 0.01	0.28	0.21	+0.07
GSPO-REPO-D	0.44 ± 0.03	0.48	0.22 ± 0.02	0.24	0.17	+0.03

Table 2: Task goal completion scores for AppWorld Qwen-3-8B by training algorithm.

5 RESULTS AND DISCUSSION

5.1 DISTINCT ALGORITHMS SHOW VARIABLE ENTROPY DYNAMICS

Results for AppWorld are in Tabs. 1 and 2 and Figs. 2 and 3. Results for AIME are in Tab. 3. We narrate these results bellow. The conclusions presented hold across all model and dataset combinations.

PPO-like algorithms deplete entropy faster than strictly on-policy. GRPO and LOOP reduce entropy by nearly 90% over training. While RLOO loses considerably less.

Clipping modifications can protect entropy. Following the intuition provided in §3, both DAPO and particularly GSPO retain more entropy than GRPO or LOOP. DAPO and GSPO are methods that have been empirically found to perform better. Here we show that one possible explanation for the improved performance of these methods is that they preserve entropy.

REPO is most effective at entropy preservation. The REPO-D and REPO-R variants, built on top of DAPO and GSPO, consistently yield no loss or even gains in entropy over training. This confirms the effectiveness of our proposed regularization and the understanding it builds upon.

REPO outperforms an entropy bonus. An entropy bonus in reinforcement learning is an additional term in the optimization objective that increase entropy directly (Williams, 1992; Mnih et al., 2016). We compare DAPO with an entropy bonus³ with REPO. While an entropy bonus aids DAPO, it is worse than REPO and uses more memory⁴ (Fig. 3).

³We use the same adaptive algorithm as REPO-D to set the coefficient β .

⁴Computing an entropy bonus requires storing the logits in GPU memory whereas computing the logprobability for the select token does not with CCE (Wijmans et al., 2025). Modifying CCE is non-trivial.

Model	Algorithm	AIME24@1	AIME24@64	AIME25@1	AIME25@64	Combined	\mathcal{H}	$\Delta \mathcal{H}$
8B	RLOO	0.62	0.79	0.44	0.67	0.63	0.280	-0.020
8B	GRPO	0.59	0.80	0.41	0.53	0.58	0.107	-0.193
8B	LOOP	0.58	0.81	0.41	0.60	0.60	0.074	-0.226
8B	DAPO	0.63	0.80	0.42	0.63	0.62	0.239	-0.061
8B	GSPO	0.63	0.82	0.43	0.67	0.64	0.193	-0.107
8B	REPO-R	0.63	0.78	0.47	0.63	0.63	0.397	+0.097
8B	REPO-D	0.63	0.83	0.42	0.63	0.63	0.353	+0.053
8B	GSPO-REPO-R	0.64	0.81	0.40	0.62	0.62	0.336	+0.036
8B	GSPO-REPO-D	0.64	0.80	0.46	0.66	0.64	0.332	+0.032
32B	RLOO	0.68	0.88	0.50	0.66	0.68	0.145	-0.235
32B	GRPO	0.64	0.82	0.48	0.64	0.65	0.047	-0.333
32B	LOOP	0.67	0.80	0.50	0.68	0.66	0.033	-0.347
32B	DAPO	0.65	0.87	0.55	0.68	0.69	0.319	-0.061
32B	GSPO	0.63	0.83	0.47	0.70	0.66	0.297	-0.083
32B	REPO-R	0.68	0.88	0.47	0.68	0.68	0.469	+0.089
32B	REPO-D	0.64	0.84	0.50	0.72	0.68	0.442	+0.062
32B	GSPO-REPO-R	0.68	0.87	0.50	0.72	0.69	0.422	+0.042
32B	GSPO-REPO-D	0.70	0.86	0.47	0.66	0.67	0.343	-0.037

Table 3: AIME results by parameter count and training algorithm.

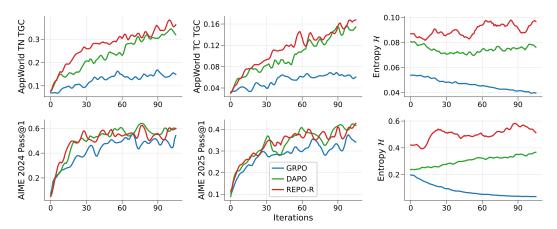


Figure 4: Sequential learning experiment. Top row: We use an AIME-trained model for GRPO, DAPO, REPO-R, and continue training the model on AppWorld. The left and middle plots show Task Goal Completion (TGC) on the normal (TN) and challenging (TC) test sets. A collapsed model (GRPO) does significantly worse than one in which entropy is preserved. Bottom row: We use an AppWorld-trained model and continue training on AIME. The same trends hold. All curves reflect the mean across multiple seeds.

5.2 Entropy preservation and downstream performance

We evaluate the effect of entropy preservation on downstream performance. See Fig. 1 for a preview of these results and App. C for a full breakdown and analysis. We find that methods that preserve per-token entropy over the course of training tend to yield higher final test accuracy than those that don't. This is also captured in the cumulative entropy over training. Methods (and checkpoints) that maintain a higher cumulative entropy over training yield a higher final test accuracy. These trends are stronger on AppWorld than AIME. We hypothesize that this is due to the following: The Qwen-3 family of models is already heavily optimized for AIME, and so this optimization may have primarily involved sharpening around existing solutions. AppWorld, on the other hand, requires discovering new capabilities (and thus requires more entropy to explore).

5.3 Entropy preservation assists sequential training

Entropy is critical for online reinforcement learning. In this regime, a policy must generate trajectories that yield different returns to collect non-zero advantage samples. If entropy is exhausted, then

learning stops. Here we ask whether entropy collapse hinders the ability of a trained policy to be re-trained in a novel environment.

We train Qwen-3-8B first on either the AIME or AppWorld using the same settings as §4. We then take the best checkpoint as the starting point for training on the opposing environment. Fig. 4 shows that policies trained with GRPO in one environment have low entropy even once transferred to the other. This results in a lower peak performance during re-training. On the other hand, DAPO, and especially REPO, start re-training with ample entropy and retain more entropy over the course of training. This results in performance comparable to starting training from Qwen-3-8B, thereby demonstrating that collapsing entropy harms re-training.

6 RELATED WORK

 Reinforcement learning has emerged as the dominant paradigm for aligning pre-trained language models (Ziegler et al., 2019; Stiennon et al., 2020; Ouyang et al., 2022). This approach has been successfully scaled in environments yielding verifiable rewards such as programming and mathematics (Jaech et al., 2024; Lambert et al., 2024; Comanici et al., 2025; Guo et al., 2025; Team et al., 2025).

Empirically, training in this setting has typically been viewed as sharpening the base policy around existing solutions rather than yielding new ones (Gandhi et al., 2025; Liu et al., 2025b; Yue et al., 2025; Zhao et al., 2025). A good pre-trained base policy starts off already calibrated to many reasonable reward functions, and post-training can be viewed as tempering this distribution (Kadavath et al., 2022; Cui et al., 2025). In fact, several works directly exploit this calibration to drive accuracy improvements via unsupervised post-training. Agarwal et al. (2024) simply minimize entropy, Prasad et al. (2024); Zhang et al. (2025); Zuo et al. (2025) align to the model's majority vote distribution, Wang et al. (2025) get by with a single labeled sample, and Shao et al. (2025) even use random rewards. All of these works can be explained by simply allowing policy gradient to sharpen an already calibrated base policy. While this type of approach can help pass@1, it harms pass@k (Shao et al., 2024; Dang et al., 2025; Yue et al., 2025).

Some works protect against this pathological entropy collapse using modified policy gradient objectives. He et al. (2025) add auxiliary rewards to solutions as a function of their probability rank within a batch. Yu et al. (2025) introduce wider PPO clipping to encourage stronger reinforcement of low probability correct actions. Zheng et al. (2025) propose sequence-level clipping more independent of individual action probabilities. Chen et al. (2025b) reformulate online policy gradient to optimize pass@k as opposed to pass@l. Most similarly to our work, Cui et al. (2025) derive theoretical results regarding the covariance between advantages and probabilities mediating entropy collapse and then identify the individual tokens most responsible for sharpening and detach their gradients.

Other works impose a $D_{\rm KL}$ penalty during training as an approach for preserving the base policy (e.g., Ziegler et al., 2019; Guo et al., 2025, etc.). However, it has been shown that such an approach limits how much the policy can learn (Korbak et al., 2022; Yang et al., 2024; Wu & Choi, 2025). For this reason, Chen et al. (2025a); Yu et al. (2025) remove the $D_{\rm KL}$ penalty, (Vassoyan et al., 2025) ignore it for a subset of tokens, and (Liu et al., 2025a) iteratively reset the reference policy.

7 Conclusion

In this work, we begin with a theoretical explanation for entropy collapse under policy gradient. We show that this process accelerates under PPO relative to strict on-policy, and how recent policy gradient variants like GSPO or DAPO, implicitly prevent this collapse. We then propose REPO, a novel approach to policy gradient optimization that uses an adaptive controller to stabilize entropy dynamics online. We provide empirical evidence for REPO's effectiveness, training in challenging environments and evaluating on AppWorld, AIME 2024, and AIME 2025. In addition to strong benchmark performance, REPO-trained models yield final policies that have retained their entropy, which we demonstrate enables sequential learning of trained checkpoints in new environments. Overall, we highlight the importance of entropy—and the corresponding online exploration capabilities—for effective policy optimization.

ETHICS STATEMENT

This paper investigates the properties of policy gradient algorithms for language model reasoning, specifically focusing on the tendency for entropy collapse during training. Our research is primarily theoretical and analytical, involving mathematical analysis and algorithm development. Our work aims to improve entropy during reinforcement learning, which can lead to better exploration and wider diversity in generated outputs. We acknowledge the potential for misuse of advanced language models, including the generation of biased, harmful, or misleading content. We believe that responsible research practices, including transparency in model limitations and potential societal impacts, are crucial for mitigating these risks, and we hope that our research contributes to the development of more robust, creative, and beneficial language models.

REPRODUCIBILITY STATEMENT

Complete proofs for all theoretical claims, along with experimental details and hyperparameters, are included in the appendix. All data points presented in this work are the result of multiple repetitions of each experiment using independent random seeds.

USE OF LARGE LANGUAGE MODELS FOR WRITING

We acknowledge the use of large language models to assist with typographical corrections, phrasing, and self-review aimed at improving the clarity and structure of this manuscript.

REFERENCES

- Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos Garea, Matthieu Geist, and Olivier Bachem. On-policy distillation of language models: Learning from self-generated mistakes. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=3zKtaqxLhW. (Cited on p. 9)
- Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting REINFORCE-style optimization for learning from human feedback in LLMs. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12248–12267, 2024. URL https://aclanthology.org/2024.acl-long.662. (Cited on p. 3)
- Kevin Chen, Marco Cusumano-Towner, Brody Huval, Aleksei Petrenko, Jackson Hamburger, Vladlen Koltun, and Philipp Krähenbühl. Reinforcement learning for long-horizon interactive LLM agents. arXiv preprint arXiv:2502.01600, 2025a. URL https://arxiv.org/abs/2502.01600. (Cited on p. 3, 9, 21)
- Zhipeng Chen, Xiaobo Qin, Youbin Wu, Yue Ling, Qinghao Ye, Wayne Xin Zhao, and Guang Shi. Pass@k training for adaptively balancing exploration and exploitation of large reasoning models. arXiv preprint arXiv:2508.10751, 2025b. URL https://arxiv.org/pdf/2508.10751. (Cited on p. 1, 9)
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. arXiv preprint arXiv:2507.06261, 2025. URL https://arxiv.org/abs/2507.06261. (Cited on p. 1, 9)
- Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen Fan, Huayu Chen, Weize Chen, et al. The entropy mechanism of reinforcement learning for reasoning language models. *arXiv preprint arXiv:2505.22617*, 2025. URL https://arxiv.org/abs/2505.22617. (Cited on p. 4, 9)
- Xingyu Dang, Christina Baek, J Zico Kolter, and Aditi Raghunathan. Assessing diversity collapse in reasoning. In *Scaling Self-Improving Foundation Models without Human Supervision*, 2025. URL https://openreview.net/forum?id=AMiksHLjQh. (Cited on p. 1, 9)

- Benjamin Eysenbach and Sergey Levine. Maximum entropy RL (provably) solves some robust RL problems. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=PtSAD3caaA2. (Cited on p. 1)
 - Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh, Nathan Lile, and Noah D Goodman. Cognitive behaviors that enable self-improving reasoners, or, four habits of highly effective stars. *arXiv* preprint arXiv:2503.01307, 2025. URL https://arxiv.org/abs/2503.01307. (Cited on p. 9)
 - Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model evaluation harness, 07 2024. URL https://zenodo.org/records/12608602. (Cited on p. 21)
 - Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025. URL https://arxiv.org/abs/2501.12948. (Cited on p. 1, 9)
 - Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with deep energy-based policies. In *International conference on machine learning*, pp. 1352–1361. PMLR, 2017. URL https://proceedings.mlr.press/v70/haarnoja17a.html. (Cited on p. 1)
 - Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. PMLR, 2018. URL https://proceedings.mlr.press/v80/haarnoja18b. (Cited on p. 1)
 - Andre He, Daniel Fried, and Sean Welleck. Rewarding the unlikely: Lifting GRPO beyond distribution sharpening. *arXiv preprint arXiv:2506.02355*, 2025. URL https://arxiv.org/abs/2506.02355. (Cited on p. 9)
 - Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=nZeVKeeFYf9. (Cited on p. 21)
 - Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. OpenAI o1 system card. *arXiv* preprint arXiv:2412.16720, 2024. URL https://arxiv.org/abs/2412.16720. (Cited on p. 1, 9)
 - Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022. URL https://arxiv.org/abs/2207.05221. (Cited on p. 9)
 - Amirhossein Kazemnejad, Milad Aghajohari, Eva Portelance, Alessandro Sordoni, Siva Reddy, Aaron Courville, and Nicolas Le Roux. VinePPO: Unlocking RL potential for LLM reasoning through refined credit assignment. In *International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=5mJrGtXVwz. (Cited on p. 3)
 - Wouter Kool, Herke van Hoof, and Max Welling. Buy 4 reinforce samples, get a baseline for free! Deep RL Meets Structured Prediction Workshop at ICLR, 2019. URL https://openreview.net/forum?id=r1lgTGL5DE. (Cited on p. 3)
 - Tomasz Korbak, Ethan Perez, and Christopher Buckley. RL with KL penalties is better viewed as Bayesian inference. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 1083–1091, 2022. URL https://aclanthology.org/2022.findings-emnlp. 77. (Cited on p. 9)

- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with PagedAttention. In *Proceedings of the 29th symposium on operating systems principles*, pp. 611–626, 2023. URL https://dl.acm.org/doi/abs/10.1145/3600006. 3613165. (Cited on p. 21)
 - Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024. URL https://arxiv.org/abs/2411.15124. (Cited on p. 9)
 - Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lample, and Stanislas Polu. NuminaMath-1.5, 2024. URL https://huggingface.co/datasets/AI-MO/NuminaMath-1.5. (Cited on p. 6)
 - Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong. ProRL: Prolonged reinforcement learning expands reasoning boundaries in large language models. *arXiv preprint arXiv:2505.24864*, 2025a. URL https://arxiv.org/abs/2505.24864. (Cited on p. 9)
 - Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding R1-Zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*, 2025b. URL https://arxiv.org/abs/2503.20783. (Cited on p. 3, 9)
 - Aaron Meurer, Christopher P Smith, Mateusz Paprocki, Ondřej Čertík, Sergey B Kirpichev, Matthew Rocklin, AMiT Kumar, Sergiu Ivanov, Jason K Moore, Sartaj Singh, et al. SymPy: symbolic computing in Python. *PeerJ Computer Science*, 3:e103, 2017. URL https://peerj.com/articles/cs-103/. (Cited on p. 21)
 - Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pp. 1928–1937. PmLR, 2016. URL https://proceedings.mlr.press/v48/mniha16.html. (Cited on p. 7)
 - Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35:27730-27744, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/hash/blefde53be364a73914f58805a001731-Abstract-Conference.html. (Cited on p. 9)
 - Archiki Prasad, Weizhe Yuan, Richard Yuanzhe Pang, Jing Xu, Maryam Fazel-Zarandi, Mohit Bansal, Sainbayar Sukhbaatar, Jason E Weston, and Jane Yu. Self-consistency preference optimization. In Forty-second International Conference on Machine Learning, 2024. URL https://openreview.net/forum?id=94G4eL3RWi. (Cited on p. 9)
 - John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pp. 1889–1897. PMLR, 2015. URL https://proceedings.mlr.press/v37/schulman15.html. (Cited on p. 3)
 - John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. URL https://arxiv.org/abs/1707.06347. (Cited on p. 3, 5)
 - Rulin Shao, Shuyue Stella Li, Rui Xin, Scott Geng, Yiping Wang, Sewoong Oh, Simon Shaolei Du, Nathan Lambert, Sewon Min, Ranjay Krishna, et al. Spurious rewards: Rethinking training signals in RLVR. arXiv preprint arXiv:2506.10947, 2025. URL https://arxiv.org/abs/2506.10947. (Cited on p. 9)

- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. DeepSeekMath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024. URL https://arxiv.org/abs/2402.03300. (Cited on p. 1, 3, 9)
 - Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in neural information processing systems*, 33:3008–3021, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/1f89885d556929e98d3ef9b86448f951-Abstract.html. (Cited on p. 9)
 - Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998. (Cited on p. 1)
 - Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with LLMs. *arXiv preprint arXiv:2501.12599*, 2025. URL https://arxiv.org/abs/2501.12599. (Cited on p. 9)
 - Sebastian B Thrun. *Efficient exploration in reinforcement learning*. Carnegie Mellon University, 1992. (Cited on p. 1)
 - Harsh Trivedi, Tushar Khot, Mareike Hartmann, Ruskin Manku, Vinty Dong, Edward Li, Shashank Gupta, Ashish Sabharwal, and Niranjan Balasubramanian. AppWorld: A controllable world of apps and people for benchmarking interactive coding agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 16022–16076, 2024. URL https://aclanthology.org/2024.acl-long.850/. (Cited on p. 5)
 - Jean Vassoyan, Nathanaël Beau, and Roman Plaud. Ignore the KL penalty! boosting exploration on critical tokens to enhance RL fine-tuning. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pp. 6108–6118, 2025. URL https://aclanthology.org/2025.findings-naacl.340. (Cited on p. 9)
 - Yiping Wang, Qing Yang, Zhiyuan Zeng, Liliang Ren, Liyuan Liu, Baolin Peng, Hao Cheng, Xuehai He, Kuan Wang, Jianfeng Gao, et al. Reinforcement learning for reasoning in large language models with one training example. arXiv preprint arXiv:2504.20571, 2025. URL https://arxiv.org/abs/2504.20571. (Cited on p. 9)
 - Erik Wijmans, Abhishek Kadian, Ari Morcos, Stefan Lee, Irfan Essa, Devi Parikh, Manolis Savva, and Dhruv Batra. DD-PPO: Learning near-perfect pointgoal navigators from 2.5 billion frames. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=H1qX8C4YPr. (Cited on p. 21)
 - Erik Wijmans, Brody Huval, Alexander Hertzberg, Vladlen Koltun, and Philipp Kraehenbuehl. Cut your losses in large-vocabulary language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=E4Fk3YuG56. (Cited on p. 7, 21)
 - Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992. URL https://link.springer.com/article/10.1007/bf00992696. (Cited on p. 3, 7, 16)
 - Fang Wu and Yejin Choi. On the limits of RLVR: Support, entropy, and the illusion of reasoning. In 2nd AI for Math Workshop@ ICML 2025, 2025. URL https://openreview.net/forum?id=KXtLWJAzgh. (Cited on p. 9)
 - An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. arXiv preprint arXiv:2505.09388, 2025. URL https://arxiv.org/abs/2505.09388. (Cited on p. 5)

- Joy Qiping Yang, Salman Salamatian, Ziteng Sun, Ananda Theertha Suresh, and Ahmad Beirami. Asymptotics of language model alignment. In 2024 IEEE International Symposium on Information Theory (ISIT), pp. 2027–2032. IEEE, 2024. URL https://ieeexplore.ieee.org/abstract/document/10619456. (Cited on p. 9)
 - Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. DAPO: An open-source LLM reinforcement learning system at scale. arXiv preprint arXiv:2503.14476, 2025. URL https://arxiv.org/abs/2503.14476. (Cited on p. 1, 4, 9)
 - Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Shiji Song, and Gao Huang. Does reinforcement learning really incentivize reasoning capacity in LLMs beyond the base model? *arXiv preprint arXiv:2504.13837*, 2025. URL https://arxiv.org/abs/2504.13837. (Cited on p. 1, 9)
 - Qingyang Zhang, Haitao Wu, Changqing Zhang, Peilin Zhao, and Yatao Bian. Right question is already half the answer: Fully unsupervised LLM reasoning incentivization. *arXiv* preprint arXiv:2504.05812, 2025. URL https://arxiv.org/abs/2504.05812. (Cited on p. 9)
 - Ruisi Zhang, Tianyu Liu, Will Feng, Andrew Gu, Sanket Purandare, Wanchao Liang, and Francisco Massa. SimpleFSDP: Simpler fully sharded data parallel with torch. compile. *arXiv* preprint *arXiv*:2411.00284, 2024. URL https://arxiv.org/abs/2411.00284. (Cited on p. 21)
 - Rosie Zhao, Alexandru Meterez, Sham Kakade, Cengiz Pehlevan, Samy Jelassi, and Eran Malach. Echo chamber: RL post-training amplifies behaviors learned in pretraining. *arXiv preprint* arXiv:2504.07912, 2025. URL https://arxiv.org/abs/2504.07912. (Cited on p. 9)
 - Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, et al. Group sequence policy optimization. *arXiv preprint arXiv:2507.18071*, 2025. URL https://arxiv.org/pdf/2507.18071. (Cited on p. 1, 3, 5, 9)
 - Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy inverse reinforcement learning. In *AAAI*, volume 8, pp. 1433–1438. Chicago, IL, USA, 2008. URL https://dl.acm.org/doi/10.5555/1620270.1620297. (Cited on p. 1)
 - Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv* preprint arXiv:1909.08593, 2019. URL https://arxiv.org/abs/1909.08593. (Cited on p. 9)
 - Yuxin Zuo, Kaiyan Zhang, Li Sheng, Shang Qu, Ganqu Cui, Xuekai Zhu, Haozhan Li, Yuchen Zhang, Xinwei Long, Ermo Hua, et al. TTRL: Test-time reinforcement learning. *arXiv* preprint arXiv:2504.16084, 2025. URL https://arxiv.org/abs/2504.16084. (Cited on p. 9)

A PROOFS & DERIVATIONS

A.1 Broadly used Lemmas

Lemma 1. The expected score function of policy π_{θ} at some state s is:

$$\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s})} \left[\nabla_{\theta} \log \pi_{\theta}(a \mid \mathbf{s}) \right] = 0$$

Proof.

$$\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s})} \left[\nabla_{\theta} \log \pi_{\theta}(a \mid \mathbf{s}) \right] = \sum_{a} \pi_{\theta}(a \mid \mathbf{s}) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \mathbf{s})$$

$$= \sum_{a} \nabla_{\theta} \pi_{\theta}(a \mid \mathbf{s})$$

$$= \nabla_{\theta} \sum_{a} \pi_{\theta}(a \mid \mathbf{s})$$

$$= \nabla_{\theta}(1)$$

$$= 0$$

Lemma 2. The gradient of a sample estimate $\mathbb{E}_{x \sim P_{\theta}}[f_{\theta}(x)]$ of function f_{θ} over distribution P_{θ} is:

$$\nabla_{\theta} \mathbb{E}_{x \sim P_{\theta}} \left[f_{\theta}(x) \right] = \mathbb{E}_{x \sim P_{\theta}} \left[\nabla_{\theta} f_{\theta}(x) + f_{\theta}(x) \cdot \nabla_{\theta} \log P_{\theta}(x) \right]$$

Proof.

$$\nabla_{\theta} \mathbb{E}_{x \sim P_{\theta}} [f_{\theta}(x)] = \sum_{x} \nabla_{\theta} (P_{\theta}(x) \cdot f_{\theta}(x))$$

$$= \sum_{x} \left(P_{\theta}(x) \cdot \nabla_{\theta} f_{\theta}(x) + f_{\theta}(x) \cdot \underbrace{\nabla_{\theta} P_{\theta}(x)}_{P_{\theta}(x) \nabla_{\theta} \log P_{\theta}(x)} \right)$$

$$= \sum_{x} P_{\theta}(x) (\nabla_{\theta} f_{\theta}(x) + f_{\theta}(x) \cdot \nabla_{\theta} \log P_{\theta}(x))$$

$$= \mathbb{E}_{x \sim P_{\theta}} [\nabla_{\theta} f_{\theta}(x) + f_{\theta}(x) \cdot \nabla_{\theta} \log P_{\theta}(x)]$$

Lemma 3. The gradient of a sample estimate $\mathbb{E}_{x \sim P_{\theta}}[f_{\theta}(x)]$ of function f_{θ} over distribution P_{θ} can be unbiasedly mean centered such that:

$$\nabla_{\theta} \mathbb{E}_{x \sim P_{\theta}} \left[f_{\theta}(x) - \mathbb{E}_{x \sim P_{\theta}} \left[f_{\theta}(x) \right] \right] = \nabla_{\theta} \mathbb{E}_{x \sim P_{\theta}} \left[f_{\theta}(x) \right]$$

Proof.

$$\nabla_{\theta} \mathbb{E}_{x \sim P_{\theta}} \left[f_{\theta}(x) - \mathbb{E}_{x \sim P_{\theta}} \left[f_{\theta}(x) \right] \right] = \mathbb{E}_{x \sim P_{\theta}} \left[\left(f_{\theta}(x) - \mathbb{E}_{x \sim P_{\theta}} \left[f_{\theta}(x) \right] \right) \cdot \nabla_{\theta} \log P_{\theta}(x) \right]$$

$$= \mathbb{E}_{x \sim P_{\theta}} \left[f_{\theta}(x) \cdot \nabla_{\theta} \log P_{\theta}(x) \right] - \mathbb{E}_{x \sim P_{\theta}} \left[\mathbb{E}_{x \sim P_{\theta}} \left[f_{\theta}(x) \right] \cdot \nabla_{\theta} \log P_{\theta}(x) \right]$$

$$= \mathbb{E}_{x \sim P_{\theta}} \left[f_{\theta}(x) \cdot \nabla_{\theta} \log P_{\theta}(x) \right] - \mathbb{E}_{x \sim P_{\theta}} \left[f_{\theta}(x) \right] \cdot \underbrace{\mathbb{E}_{x \sim P_{\theta}} \left[\nabla_{\theta} \log P_{\theta}(x) \right]}_{0}$$

$$= \mathbb{E}_{x \sim P_{\theta}} \left[f_{\theta}(x) \cdot \nabla_{\theta} \log P_{\theta}(x) \right]$$

$$= \nabla_{\theta} \mathbb{E}_{x \sim P_{\theta}} \left[f_{\theta}(x) \right]$$

Lemma 4. The gradient of the policy entropy at some state s is:

$$\nabla_{\theta} \mathcal{H}_{\pi_{\theta}}(\cdot \mid \mathbf{s}) = -\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s})} \left[(\log \pi_{\theta}(a \mid \mathbf{s}) + b) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \mathbf{s}) \right]$$

for an arbitrary baseline b.

Proof.

$$\nabla_{\theta} \mathcal{H}_{\pi_{\theta}}(\cdot \mid \mathbf{s}) = -\nabla_{\theta} \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s})} \left[\log \pi_{\theta}(a \mid \mathbf{s}) \right]$$

$$= -\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s})} \left[\log \pi_{\theta}(a \mid \mathbf{s}) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \mathbf{s}) \right] - \underbrace{\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s})} \left[\nabla_{\theta} \log \pi_{\theta}(a \mid \mathbf{s}) \right]}_{0}$$

$$= -\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s})} \left[\log \pi_{\theta}(a \mid \mathbf{s}) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \mathbf{s}) \right]$$

$$= -\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s})} \left[\log \pi_{\theta}(a \mid \mathbf{s}) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \mathbf{s}) \right] + b \underbrace{\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s})} \left[\nabla_{\theta} \log \pi_{\theta}(a \mid \mathbf{s}) \right]}_{0}$$

$$= -\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s})} \left[(\log \pi_{\theta}(a \mid \mathbf{s}) + b) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \mathbf{s}) \right]$$

Lemma 5. The gradient of MDP objective \mathcal{J}_{MDP} at some state s is:

$$\nabla_{\theta} \mathcal{J}_{\text{MDP}}(\boldsymbol{s}) = \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[R(\boldsymbol{s}, a) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \boldsymbol{s}) \right]$$

Proof. Largely following (Williams, 1992) and Lemma 2,

$$\nabla_{\theta} \mathcal{J}_{\text{MDP}}(\boldsymbol{s}) = \nabla_{\theta} \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} [R(\boldsymbol{s}, a)]$$

$$= \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} [R(\boldsymbol{s}, a) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \boldsymbol{s})] + \underbrace{\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} [\nabla_{\theta} R(\boldsymbol{s}, a)]}_{0}$$

$$= \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} [R(\boldsymbol{s}, a) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \boldsymbol{s})]$$

Lemma 6. The advantage function $A(s, a) \stackrel{\text{def}}{=} R(s, a) - b$, with baseline independent of a, yields an unbiased policy gradient.

Proof.

$$\nabla_{\theta} \mathcal{J}_{ADV}(\boldsymbol{s}) = \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[A(\boldsymbol{s}, a) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \boldsymbol{s}) \right]$$

$$= \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[(R(\boldsymbol{s}, a) - b) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \boldsymbol{s}) \right]$$

$$= \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[R(\boldsymbol{s}, a) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \boldsymbol{s}) \right] - \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[b \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \boldsymbol{s}) \right]$$

$$= \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[R(\boldsymbol{s}, a) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \boldsymbol{s}) \right] - b \cdot \underbrace{\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[\nabla_{\theta} \log \pi_{\theta}(a \mid \boldsymbol{s}) \right]}_{0}$$

$$= \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[R(\boldsymbol{s}, a) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \boldsymbol{s}) \right]$$

$$= \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[R(\boldsymbol{s}, a) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \boldsymbol{s}) \right]$$

$$= \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[R(\boldsymbol{s}, a) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid \boldsymbol{s}) \right]$$

Lemma 7. The expected advantage function, with baseline $V(s) \stackrel{\text{def}}{=} \mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)}[R(s, a)]$, at some state s is:

$$\mathbb{E}_{a \sim \pi_{\theta}(\cdot | \boldsymbol{s})}[A(\boldsymbol{s}, a)] = 0$$

Proof.

$$\mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)}[A(s, a)] = \mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)}[R(s, a) - V(s)]$$

$$= \mathbb{E}_{a \sim \pi_{\theta}(\cdot | s)}[R(s, a)] - V(s)$$

$$= V(s) - V(s)$$

$$= 0$$

A.2 ENTROPY DYNAMICS UNDER POLICY GRADIENT

Theorem 1. Given a policy gradient update $\widehat{\theta} := \theta + \alpha \cdot \nabla_{\theta} \mathcal{J}_{MDP}(s)$, the expected change in entropy is approximately:

$$\Delta \mathcal{H}_{\pi_{\theta}}(\cdot \mid \boldsymbol{s}) \approx -\alpha \cdot \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s}), a' \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[A(\boldsymbol{s}, a) \cdot \left(\log \pi_{\theta}(a' \mid \boldsymbol{s}) - \overline{\log \pi_{\theta}(\cdot \mid \boldsymbol{s})} \right) \cdot u(\boldsymbol{s}, a)^{\top} u(\boldsymbol{s}, a') \right],$$

where $\overline{\log \pi_{\theta}(\cdot \mid s)} = \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid s)}[\log \pi_{\theta}(a \mid s)]$ and $u(s, a) \stackrel{\text{def}}{=} \nabla_{\theta} \log \pi_{\theta}(a \mid s)$ is the score function of π_{θ} evaluated at state s and action a.

Proof. Let $u(s,a) \stackrel{\text{def}}{=} \nabla_{\theta} \log \pi_{\theta}(a \mid s)$ denote the score function of policy π_{θ} evaluated at action a and state s. Let g(s) and h(s) denote the respective mean-baselined policy gradient and entropy gradient evaluated on-policy in some state s:

$$g(s) = \nabla_{\theta} \mathcal{J}_{\text{MDP}}(s) = \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid s)} \left[A(s, a) \cdot u(s, a) \right]$$
$$h(s) = \nabla_{\theta} \mathcal{H}_{\pi_{\theta}}(\cdot \mid s) = -\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid s)} \left[\left(\log \pi_{\theta}(a \mid s) - b_{s} \right) \cdot u(s, a) \right]$$

Here the entropy estimator allows for an arbitrary baseline that cancels through the parameter gradient ∇_{θ} . While the baseline does not influence the exact mathematical construction, it does influence approximations to the change in entropy. Here we chose $b_{s} = \overline{\log \pi_{\theta}(\cdot \mid s)} = E_{a \sim \pi_{\theta}(\cdot \mid s)} [\log \pi_{\theta}(a \mid s)]$, as it centers the policy and agrees with a tabular softmax approximation of the change in entropy, see Corollary 2.

Using the first-order Taylor approximation: $\mathcal{H}_{\pi_{\theta}}(\cdot \mid \mathbf{s} ; \theta + \alpha \cdot g) \approx \mathcal{H}_{\pi_{\theta}}(\cdot \mid \mathbf{s} ; \theta) + \alpha \cdot g^{\top}h$, for small learning rate α , the expected change in entropy from a policy gradient update in state \mathbf{s} is:

$$\Delta \mathcal{H}_{\pi_{\theta}}(\cdot \mid \mathbf{s}) \approx \alpha \cdot g(\mathbf{s})^{\top} h(\mathbf{s})$$

$$= -\alpha \cdot \left(\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s})} \left[A(\mathbf{s}, a) \cdot u(\mathbf{s}, a) \right] \right)^{\top} \left(\mathbb{E}_{a' \sim \pi_{\theta}(\cdot \mid \mathbf{s})} \left[\left(\log \pi_{\theta}(a' \mid \mathbf{s}) - \overline{\log \pi_{\theta}(\cdot \mid \mathbf{s})} \right) \cdot u(\mathbf{s}, a') \right] \right)$$

$$= -\alpha \cdot \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s}), a' \sim \pi_{\theta}(\cdot \mid \mathbf{s})} \left[A(\mathbf{s}, a) \cdot \left(\log \pi_{\theta}(a' \mid \mathbf{s}) - \overline{\log \pi_{\theta}(\cdot \mid \mathbf{s})} \right) \cdot u(\mathbf{s}, a)^{\top} u(\mathbf{s}, a') \right]$$

A.3 APPROXIMATE ENTROPY DYNAMICS UNDER POLICY GRADIENT

Corollary 1. Assuming $u(s, a)^{\top}u(s, a') = 0$ for all $a \neq a'$, the change in entropy is approximately:

$$\Delta \mathcal{H}_{\pi_{\theta}}(\cdot \mid \boldsymbol{s}) \propto -\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[A(\boldsymbol{s}, a) \cdot \left(\log \pi_{\theta}(a \mid \boldsymbol{s}) - \overline{\log \pi_{\theta}(\cdot \mid \boldsymbol{s})} \right) \cdot \pi_{\theta}(a \mid \boldsymbol{s}) \right]$$

Proof. Assuming the score vectors satisfy orthogonality of the off-diagonal terms such that $u(s, a)^{\top}u(s, a') = 0$ for $a \neq a'$, the double expectation can be collapsed, yielding:

$$\Delta \mathcal{H}_{\pi_{\theta}}(\cdot \mid \boldsymbol{s}) \approx -\alpha \cdot \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \left[\pi_{\theta}(a \mid \boldsymbol{s}) \cdot A(\boldsymbol{s}, a) \cdot \left(\log \pi_{\theta}(a \mid \boldsymbol{s}) - \overline{\log \pi_{\theta}(\cdot \mid \boldsymbol{s})} \right) \cdot \|u(\boldsymbol{s}, a)\|^{2} \right]$$

Assuming independence of the squared gradient norm magnitude, such that it can be treated as a constant with respect to the expectation,

$$\Delta \mathcal{H}_{\pi_{\theta}}(\cdot \mid \boldsymbol{s}) \propto -\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})}[A(\boldsymbol{s}, a) \cdot \left(\log \pi_{\theta}(a \mid \boldsymbol{s}) - \overline{\log \pi_{\theta}(\cdot \mid \boldsymbol{s})}\right) \cdot \pi_{\theta}(a \mid \boldsymbol{s})]$$

A.4 Entropy dynamics under policy gradient for tabular softmax policies

Proposition 1. For two functions f(x) and g(x) over samples $x \sim \pi_S$ of a softmax distribution $\pi_S(x) = \exp(S_x) / \sum_k \exp(S_k)$, the dot product of expected gradients is:

$$\left\langle \mathbb{E}_{x \sim \pi_S}[f(x) \cdot \nabla_S \log \pi_S(x)] \right. , \quad \mathbb{E}_{y \sim \pi_S}[g(y) \cdot \nabla_S \log \pi_S(y)] \right\rangle = \mathbb{E}_{x \sim \pi_S}[\pi_S(x) \cdot (f(x) - \bar{f}) \cdot (g(x) - \bar{g})],$$
 where $\bar{f} = \mathbb{E}_{x \sim \pi_S}[f(x)]$ and $\bar{g} = \mathbb{E}_{x \sim \pi_S}[g(x)].$

Proof. First, let's compute $\nabla_S \log \pi_S(x)$ for the softmax distribution:

$$\log \pi_S(x) = \log \frac{\exp(S_x)}{\sum_k \exp(S_k)} = S_x - \log \sum_k \exp(S_k)$$
$$\nabla_{S_z} \log \pi_S(x) = \mathbb{1}_{x=z} - \frac{\exp(S_z)}{\sum_k \exp(S_k)} = \mathbb{1}_{x=z} - \pi_S(z)$$

where $\mathbb{1}_{x=y}$ is the indicator function (1 if x=y,0 otherwise).

Now let's compute the dot product $\nabla_S \log \pi_S(x)^\top \nabla_S \log \pi_S(y)$:

$$\nabla_{S} \log \pi_{S}(x)^{\top} \nabla_{S} \log \pi_{S}(y) = \sum_{j} (\mathbb{1}_{x=j} - \pi_{S}(j)) (\mathbb{1}_{y=j} - \pi_{S}(j))$$

$$= \sum_{j} (\mathbb{1}_{x=j} \cdot \mathbb{1}_{y=j} - \mathbb{1}_{x=j} \pi_{S}(j) - \pi_{S}(j) \cdot \mathbb{1}_{y=j} + \pi_{S}(j)^{2})$$

$$= \mathbb{1}_{x=y} - \pi_{S}(x) - \pi_{S}(y) + E_{z \sim \pi_{S}}[\pi_{S}(z)]$$

Now we can compute the dot product of expected gradients:

$$\left\langle \mathbb{E}_{x \sim \pi_S}[f(x) \cdot \nabla_S \log \pi_S(x)], \mathbb{E}_{y \sim \pi_S}[g(y) \cdot \nabla_S \log \pi_S(y)] \right\rangle$$

$$= \mathbb{E}_{x \sim \pi_S, y \sim \pi_S}[f(x) \cdot g(y) \cdot \nabla_S \log \pi_S(x)^{\top} \nabla_S \log \pi_S(y)]$$

$$= \mathbb{E}_{x \sim \pi_S, y \sim \pi_S}[f(x) \cdot g(y) \cdot (\mathbb{1}_{x=y} - \pi_S(x) - \pi_S(y) + \mathbb{E}_{z \sim \pi_S}[\pi_S(z)])]$$

Let's compute each term separately:

$$\mathbb{E}_{x \sim \pi_S, y \sim \pi_S}[f(x) \cdot g(y) \cdot \mathbb{1}_{x=y}] = \mathbb{E}_{x \sim \pi_S}[\pi_S(x) \cdot f(x) \cdot g(x)]$$

$$\mathbb{E}_{x \sim \pi_S, y \sim \pi_S}[f(x) \cdot g(y) \cdot \pi_S(x)] = \mathbb{E}_{x \sim \pi_S}[f(x) \cdot \pi_S(x)] \cdot \mathbb{E}_{y \sim \pi_S}[g(y)] = \mathbb{E}_{x \sim \pi_S}[\pi_S(x) \cdot f(x)] \cdot \bar{g}$$

$$\mathbb{E}_{x \sim \pi_S, y \sim \pi_S}[f(x) \cdot g(y) \cdot \pi_S(y)] = \mathbb{E}_{x \sim \pi_S}[f(x)] \cdot \mathbb{E}_{y \sim \pi_S}[g(y) \cdot \pi_S(y)] = \bar{f}\mathbb{E}_{y \sim \pi_S}[\pi_S(y) \cdot g(y)]$$

$$\mathbb{E}_{x \sim \pi_S, y \sim \pi_S}[f(x) \cdot g(y) \cdot \mathbb{E}_{z \sim \pi_S}[\pi_S(z)]] = \mathbb{E}_{z \sim \pi_S}[\pi_S(z)] \cdot \mathbb{E}_{x \sim \pi_S}[f(x)] \cdot \mathbb{E}_{y \sim \pi_S}[g(y)] = \mathbb{E}_{x \sim \pi_S}[\pi_S(x)] \cdot \bar{f} \cdot \bar{g}$$

Therefore:

$$\begin{split} &\left\langle \mathbb{E}_{x \sim \pi_S}[f(x) \cdot \nabla_S \log \pi_S(x)], \mathbb{E}_{y \sim \pi_S}[g(y) \cdot \nabla_S \log \pi_S(y)] \right\rangle \\ &= \mathbb{E}_{x \sim \pi_S}[\pi_S(x) \cdot f(x) \cdot g(x)] - \mathbb{E}_{x \sim \pi_S}[\pi_S(x) \cdot f(x)] \cdot \bar{g} - \bar{f} \cdot \mathbb{E}_{x \sim \pi_S}[\pi_S(x) \cdot g(x)] + \mathbb{E}_{x \sim \pi_S}[\pi_S(x)] \cdot \bar{f} \cdot \bar{g} \\ &= \mathbb{E}_{x \sim \pi_S}[\pi_S(x) \cdot (f(x) \cdot g(x) - f(x) \cdot \bar{g} - \bar{f} \cdot g(x) + \bar{f} \cdot \bar{g})] \\ &= \mathbb{E}_{x \sim \pi_S}[\pi_S(x) \cdot (f(x) - \bar{f}) \cdot (g(x) - \bar{g})] \\ \text{where } \bar{f} = \mathbb{E}_{x \sim \pi_S}[f(x)] \text{ and } \bar{g} = \mathbb{E}_{x \sim \pi_S}[g(x)]. \end{split}$$

The above proposition holds for simple softmax policies, but involves a much more complex gradient term and inner product for generic transformer-based policies.

Corollary 2. Under a tabular softmax policy, a policy gradient update $\hat{\theta} := \theta + \alpha \cdot \nabla_{\theta} \mathcal{J}_{\text{MDP}}$ changes the entropy approximately:

$$\Delta \mathcal{H}_{\pi_{\theta}}(\cdot \mid \mathbf{s}) \approx -\alpha \cdot \mathbb{E}_{a \sim \pi_{S}(\cdot \mid \mathbf{s})} \left[\pi_{S}(a \mid \mathbf{s}) \cdot \left(\log \pi_{S}(a \mid \mathbf{s}) - \overline{\log \pi_{S}} \right) \cdot A(\mathbf{s}, a) \right]$$

where $\overline{\log \pi_{\theta}} = \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \mathbf{s})} [\log \pi_{\theta}(a \mid \mathbf{s})].$

Proof. Let g(s) and h(s) denote the respective policy gradient and entropy gradient evaluated on-policy in some state s:

$$g(s) = \nabla_{\theta} \mathcal{J}_{\text{MDP}}(s) = \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid s)} \left[A(s, a) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid s) \right]$$
$$h(s) = \nabla_{\theta} \mathcal{H}_{\pi_{\theta}}(\cdot \mid s) = -\mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid s)} \left[\log \pi_{\theta}(a \mid s) \cdot \nabla_{\theta} \log \pi_{\theta}(a \mid s) \right]$$

 Using the first-order Taylor approximation: $\mathcal{H}_{\pi_{\theta}}(\cdot \mid s ; \theta + \alpha \cdot g) \approx \mathcal{H}_{\pi_{\theta}}(\cdot \mid s ; \theta) + \alpha \cdot g^{\top}h$, for small learning rate α , the expected change in entropy from a policy gradient update in state s is:

$$\Delta \mathcal{H}_{\pi_{\theta}}(\cdot \mid \mathbf{s}) \approx \alpha \cdot g(\mathbf{s})^{\top} h(\mathbf{s})$$

$$= -\alpha \cdot \mathbb{E}_{a \sim \pi_{S}(\cdot \mid \mathbf{s})} \left[\pi_{S}(a \mid \mathbf{s}) \cdot \left(\log \pi_{S}(a \mid \mathbf{s}) - \overline{\log \pi_{S}} \right) \cdot \left(A(\mathbf{s}, a) - \overline{A} \right) \right]$$

$$= -\alpha \cdot \mathbb{E}_{a \sim \pi_{S}(\cdot \mid \mathbf{s})} \left[\pi_{S}(a \mid \mathbf{s}) \cdot \left(\log \pi_{S}(a \mid \mathbf{s}) - \overline{\log \pi_{S}} \right) \cdot A(\mathbf{s}, a) \right]$$

Here, the second line follows Prop. 1 and the third line follows Lemma 7.

A.5 ENTROPY DYNAMICS UNDER CLIPPED PPO

Proposition 2. Given two distributions $\pi(x)$ and $\phi(x)$ with constraint $\frac{\pi(x)}{\phi(x)} \leq 1 + \varepsilon$ for all x, their relative entropy is bound by

$$\mathcal{H}(\pi) \le (1+\varepsilon)\mathcal{H}(\phi)$$

Proof. Let's parametrize $\pi(x) = \beta_x \phi(x)$ with $\beta_x \ge 0$ and compute its probability

$$\mathcal{H}(\pi) = -\mathbb{E}_{x \sim \pi} \left[\log \pi(a) \right]$$

$$= -\mathbb{E}_{x \sim \pi} \left[\log \phi(a) \right] - \mathbb{E}_{x \sim \pi} \left[\log \beta_x \right]$$

$$= -\mathbb{E}_{x \sim \pi} \left[\log \phi(a) \right] - \underbrace{\mathbb{E}_{x \sim \pi} \left[\log \frac{\pi(x)}{\phi(x)} \right]}_{D_{\text{KL}}(\pi \| \phi) \geq 0}$$

$$\leq -\mathbb{E}_{x \sim \pi} \left[\log \phi(a) \right]$$

$$= -\mathbb{E}_{x \sim \phi} \left[\frac{\pi(a)}{\phi(a)} \log \phi(a) \right]$$

$$= \mathbb{E}_{x \sim \phi} \left[\beta_x \cdot -\log \phi(a) \right]$$

$$\leq \mathbb{E}_{x \sim \phi} \left[(1 + \varepsilon) \cdot -\log \phi(a) \right]$$

$$= (1 + \varepsilon) \cdot \mathcal{H}(\phi)$$

The second-last line uses $-\log \phi(a) \ge 0$ and $\beta_x \le (1+\varepsilon)$ by definition, hence $\beta_x \cdot -\log \phi(a) \le (1+\varepsilon) \cdot -\log \phi(a)$.

Theorem 2. Proximal Policy Optimization (PPO) bounds the entropy $\mathcal{H}_{\pi_{\theta}^{new}}(\cdot \mid s)$ of the updated policy by the original policy entropy $\mathcal{H}_{\pi_{\theta}^{old}}(\cdot \mid s)$ such that:

$$(1 - \epsilon_{low}) \cdot \mathcal{H}_{\pi_o old}(\cdot \mid s) \leq \mathcal{H}_{\pi_o new}(\cdot \mid s) \leq (1 + \epsilon_{high}) \cdot \mathcal{H}_{\pi_o old}(\cdot \mid s)$$

Proof. Applying Prop. 2 to $\frac{\pi_{\theta}^{\text{ new}}}{\pi_{\theta}^{\text{ old}}} \leq 1 + \epsilon_{\text{high}}$ yields the upper bound

$$\mathcal{H}_{\pi_{\theta}^{\text{new}}}(\cdot \mid \boldsymbol{s}) \leq (1 + \epsilon_{\text{high}}) \cdot \mathcal{H}_{\pi_{\theta}^{\text{old}}}(\cdot \mid \boldsymbol{s}).$$

Applying Prop. 2 to $1 - \epsilon_{\text{low}} \le \frac{\pi_{\theta}^{\text{new}}}{\pi_{\theta}^{\text{old}}}$ (eqivalently $\frac{\pi_{\theta}^{\text{old}}}{\pi_{\theta}^{\text{new}}} \le \frac{1}{1 - \epsilon_{\text{low}}}$ yields the lower bound

$$\mathcal{H}_{\pi_{\boldsymbol{\theta}}^{\text{old}}}(\cdot \mid \boldsymbol{s}) \leq \frac{1}{1 - \epsilon_{\text{low}}} \cdot \mathcal{H}_{\pi_{\boldsymbol{\theta}}^{\text{new}}}(\cdot \mid \boldsymbol{s})$$

or equivalently

$$(1 - \epsilon_{\mathrm{low}}) \cdot \mathcal{H}_{\pi_{\theta}^{\mathrm{old}}}(\cdot \mid \boldsymbol{s}) \leq \mathcal{H}_{\pi_{\theta}^{\mathrm{new}}}(\cdot \mid \boldsymbol{s}).$$

A.6 Entropy change with A_{REPO} advantage function **Proposition 3.** Let $x_c(s, a) \stackrel{\text{def}}{=} \log \pi_{\theta}(a \mid s) - \overline{\log \pi_{\theta}(\cdot \mid s)},$ and define the REPO advantage $A_{\text{REPO}}(\boldsymbol{s}, a) \stackrel{\text{def}}{=} A(\boldsymbol{s}, a) - \beta_{\boldsymbol{s}} x_{c}(\boldsymbol{s}, a).$ Under the setting of Thm. 1, the first-order entropy change induced by a policy-gradient step using A_{REPO} is $\Delta \mathcal{H}_{\pi_{\theta}}^{\text{REPO}}(\cdot \mid \boldsymbol{s}) \approx \Delta \mathcal{H}_{\pi_{\theta}}(\cdot \mid \boldsymbol{s}) + \beta_{\boldsymbol{s}} \alpha \| \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})}[x_c(\boldsymbol{s}, a) \ u(\boldsymbol{s}, a)] \|^2,$ *Proof.* By Thm. 1, with the baseline choice $b_s = \overline{\log \pi_{\theta}(\cdot \mid s)}$, $h(\mathbf{s}) = -\mathbb{E}_{a \sim \pi_{\theta}(\cdot|\mathbf{s})} [x_c(\mathbf{s}, a) \ u(\mathbf{s}, a)].$ Using A_{REPO} , the policy–gradient direction becomes $g_{\text{REPO}}(\boldsymbol{s}) = \mathbb{E}_{a \sim \pi_{\theta}(\cdot \mid \boldsymbol{s})} \big[A_{\text{REPO}}(\boldsymbol{s}, a) \ u(\boldsymbol{s}, a) \big] = \mathbb{E}[A \, u] \ - \ \beta_{\boldsymbol{s}} \, \mathbb{E}[x_c \, u] = g(\boldsymbol{s}) \ - \ \beta_{\boldsymbol{s}} \, \mathbb{E}[x_c \, u].$ The first-order entropy change is $\Delta \mathcal{H}_{\pi_{\boldsymbol{\theta}}}^{\text{REPO}}(\cdot \mid \boldsymbol{s}) \approx \alpha g_{\text{REPO}}(\boldsymbol{s})^{\top} h(\boldsymbol{s}) = \alpha (g(\boldsymbol{s}) - \beta_{\boldsymbol{s}} \mathbb{E}[x_c u])^{\top} h(\boldsymbol{s}).$ Substitute $h(s) = -\mathbb{E}[x_c u]$: $\Delta \mathcal{H}_{\pi_{\theta}}^{\text{REPO}}(\cdot \mid \boldsymbol{s}) = \alpha g(\boldsymbol{s})^{\top} (-\mathbb{E}[x_c \, u]) + \beta_{\boldsymbol{s}} \alpha \mathbb{E}[x_c \, u]^{\top} \mathbb{E}[x_c \, u] = \Delta \mathcal{H}_{\pi_{\theta}}(\cdot \mid \boldsymbol{s}) + \beta_{\boldsymbol{s}} \alpha \|\mathbb{E}[x_c \, u]\|^2$

B ADDITIONAL EXPERIMENT DETAILS

B.1 Environments

Interactive tool-use agent. For the AppWorld benchmark, rollouts proceed in turns (up to 30 turns during training and 50 during evaluation), in a manner akin to an interactive notebook.

During each turn, the model generations are parsed to extract any Python code blocks, potentially containing calls to AppWorld API. These are executed to retrieve information or alter the environment state. The outputs of successful API calls or the error trace of incorrect calls appear in the agent's context after each turn. Once done, the agent may mark a task as completed at which point it is assessed whether the task state was updated successfully. Failure to mark the task as complete within the turn or context limit (32K) results in a failure. Sparse outcome-based rewards in [0,1] are assigned during training as the fraction of passing unit-tests. Binary rewards in $\{0,1\}$ are used during evaluation requiring complete correctness.

Mathematical reasoning. For the AIME benchmarks, model responses are processed and scored using the Eleuther AI lm-eval-harness Minerva math parsing utilities (Gao et al., 2024). The final unnormalized answer is first identified and parsed, then the answer is normalized to remove units, formatting, etc., and finally equivalence between the model answer and reference answer is determined using Sympy (Meurer et al., 2017).

B.2 TRAINING

Experiments are executed on 3 NVIDIA H100 8-GPU nodes. One node is used for rollout generation one for learning, and one for evaluation. Rollouts are generated using two instances of vLLM (Kwon et al., 2023) servers using 4 GPUs each with tensor parallelism. Custom RL implementation based on FSDP2 (Zhang et al., 2024) is used for training. To account for any discrepancies between sampling and training subsystems, the log-probabilities of rollout tokens are recalculated on the training node to ensure accurate importance weights for backpropagation. Cut-Cross-Entropy (CCE) is used to reduce the memory footprint during training by preventing the materialization of all logits except the target (Wijmans et al., 2025). Models are fine-tuned with LoRA (rank = 16, $\alpha = 32$) on the self-attention (key, value, query, output) and MLP modules (Hu et al., 2022). We use an AdamW optimizer with a constant learning rate of 5×10^{-5} , weight-decay = 0.01, and gradient clipping with max-norm = 0.1. To speed up rollout collection, we introduce an early stopping criteria. Once at least 4/6 rollouts per task and 90% of total rollouts are collected, we immediately proceed to training to prevent bottlenecks caused by very few extra long generations (Wijmans et al., 2020; Chen et al., 2025a).

C ADDITIONAL RESULTS

C.1 GEOMETRIC INTERPRETATION OF REPO

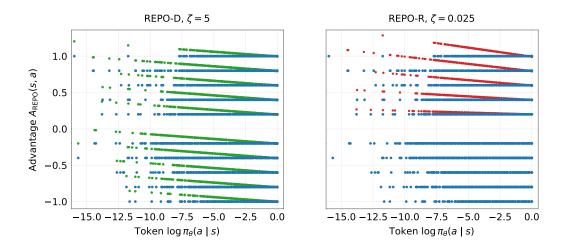


Figure 5: The REPO transformation rotates $(A, \log \pi)$ pairs, promoting low probability actions. Original unmodified advantages shown in blue, advantages A_{REPO} are shown in green for REPO-D and in red for REPO-R.

The transformation induced by each REPO algorithm can be viewed in Fig. 5. REPO-D reflects a consistent rotation across the space, boosting the advantages of actions proportional to their surprisals $(-\log \pi_\theta)$. REPO-R instead rotates only positive advantage actions, and does so proportionally, not only to the surprisal, but to the magnitude of the advantage. This strongly reinforces low-probability correct actions, especially when they yield outcomes significantly better than average for a given batch of experience.

Fig. 5 uses data from the <code>Qwen-3-8B</code> AIME experiment. The parameters of the algorithm are revealed in the structure of the data: there are 5 distinct positive and negative advantage values, corresponding to 5 unique outcomes of group-based advantage estimation (1 success / 5 failures, 2 successes / 4 failures, etc.). Groups with zero advantages are filtered out.

With the appropriate value of ζ , REPO-D transformation counteracts the covariance-like term in $\Delta \mathcal{H}$ approximation, therefore REPO-D is short for REPO-Decorrelate. REPO-R is a shorthand for REPO-Rescale, as it rescales the advantages by $1 - \zeta(\log \pi_{\theta}(a_t|s) - \overline{\log \pi_{\theta}(\cdot|s)})$

C.2 DYNAMICS OF ENTROPY AND TEST-ACCURACY DURING TRAINING

We computed the average per-token entropy at each iteration of training for all of our training runs (averaging over all tokens generated during a rollout, and averaging over all rollouts at a given iteration). We studied how this quantity evolved over the course of training for several baseline algorithms (RLOO, GRPO, LOOP, GSPO) and several variants of our REPO algorithm (REPO-R, REPO-D, GSPO-REPO-R, and GSPO-REPO-D). Figure 6 shows how the per-token entropy coevolves with the test accuracy over the course of training for each algorithm. We first observe that the REPO algorithms typically preserve much higher entropy than baselines. For challenging model-task pairs where baselines reduce policy entropy prior to achieving high test accuracy (e.g. Qwen-3-8B on all tasks, and Qwen-3-32B on AppWorld Test-Normal), REPO algorithms preserve entropy for longer and achieve higher peak test accuracy. For model-task pairs where the baselines reduce entropy late in training, after test accuracy is largely saturated, REPO achieves comparable peak test accuracy to baselines.

C.3 DEPENDENCE OF TEST ACCURACY ON CUMULATIVE ENTROPY DURING TRAINING

We hypothesized that the test accuracy at a given checkpoint is highly dependent on the *cumulative* entropy (intuitively, the time integral of the average per-token entropy) experienced over the course of previous training iterations. Figure 7 plots the test accuracy and cumulative entropy for each checkpoint of each training run of each learning algorithm studied. We observe test accuracy on a more difficult learning task (AppWorld Test-Normal) show sustained increases in test accuracy with additional cumulative entropy even late into training, whereas AIME24 and AIME25 require less cumulative entropy to achieve peak test accuracy. We quantified the dependence of test accuracy on cumulative entropy via mutual information (Table 4). We estimated (i) the mutual information between the test accuracy and the cumulative entropy and (ii) the mutual information between the test accuracy and the iteration using histograms. We found that cumulative entropy is more predictive than the iteration number. We also confirm the relatively stronger dependence of test accuracy on cumulative entropy in AppWorld Test-Normal (MI=0.858 for Qwen-3-8B and MI=0.612 for Qwen-3-32B) compared to the AIME24 and AIME25 environments (MI \approx 0.2 for both models).

	MI (Iteration)	MI (Cumulative Entropy)
Qwen3 8B - AIME25	0.182	0.205
Qwen3 32B - AIME25	0.183	0.191
Qwen3 8B - AIME24	0.131	0.170
Qwen3 32B - AIME24	0.146	0.133
Qwen3 8B - AppWorld Test-Normal	0.566	0.858
Qwen3 32B - AppWorld Test-Normal	0.507	0.612

Table 4: Quantifying the dependence of test accuracy on cumulative entropy during training.

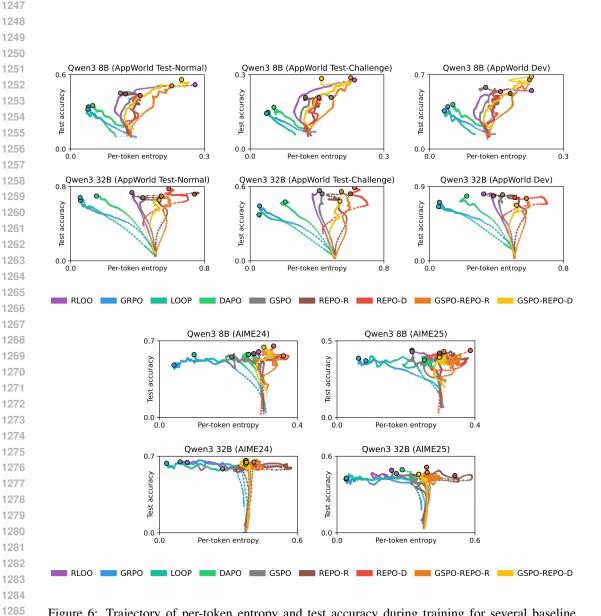


Figure 6: Trajectory of per-token entropy and test accuracy during training for several baseline algorithms several REPO algorithms in the AppWorld environment (top grid) and AIME (bottom grid). Each curve shows the average trajectory over multiple training runs with different seeds.

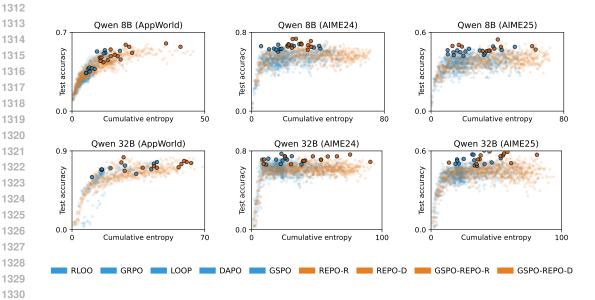


Figure 7: Accuracy on test set versus cumulative per-token entropy (sum of average per-token entropies during the training run up to that point) for all training checkpoints. Dark points show checkpoints with peak test accuracy. AppWorld test accuracy is measured on the Test-Normal set.