

Learning to Reason under Off-Policy Guidance

Jianhao Yan^{21*} Yafu Li^{1*} Zican Hu³¹ Zhi Wang³ Ganqu Cui¹ Xiaoye Qu¹
Yu Cheng^{4†} Yue Zhang^{2†}

¹ Shanghai AI Laboratory ² Westlake University ³ Nanjing University

⁴ The Chinese University of Hong Kong

Corresponding to: chengyu@cse.cuhk.edu.hk, yue.zhang@wias.org.cn

Project Page: <https://github.com/ElliottYan/LUFFY>

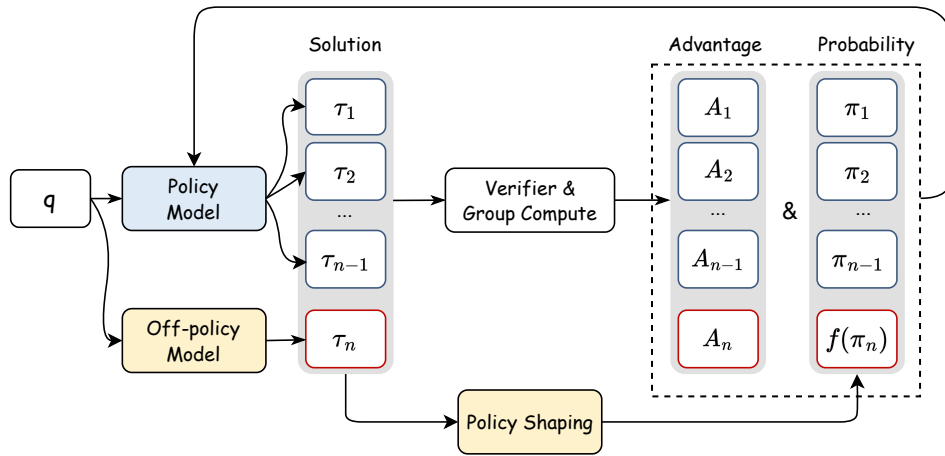


Figure 1: Overview: LUFFY integrates off-policy reasoning traces into reinforcement learning by combining them with on-policy rollouts. Policy shaping emphasizes low-probability but crucial actions, enabling a balance between imitation and exploration for more generalizable reasoning.

Abstract

Recent advances in large reasoning models (LRMs) demonstrate that sophisticated behaviors such as multi-step reasoning and self-reflection can emerge via reinforcement learning (RL) with simple rule-based rewards. However, existing *zero-RL* approaches are inherently “on-policy”, limiting learning to a model’s own outputs and failing to acquire reasoning abilities beyond its initial capabilities. We introduce **LUFFY** (Learning to reason Under oFF-policyY guidance), a framework that augments *zero-RL* with off-policy reasoning traces. LUFFY dynamically balances imitation and exploration by combining off-policy demonstrations with on-policy rollouts during training. Notably, we propose policy shaping via regularized importance sampling to avoid superficial and rigid imitation during mixed-policy training. Remarkably, LUFFY achieves an over **+7.0** average gain across six math benchmarks and an advantage of over **+6.2** points in out-of-distribution tasks. It also substantially surpasses imitation-based supervised fine-tuning (SFT), particularly in generalization. Analysis shows LUFFY not only imitates effectively but also explores beyond demonstrations, offering a scalable path to train generalizable reasoning models with off-policy guidance.

* Equal contributions. Work was done during Jianhao Yan’s internship at Shanghai AI Laboratory. Yafu Li is the Project Lead.

† Corresponding authors.

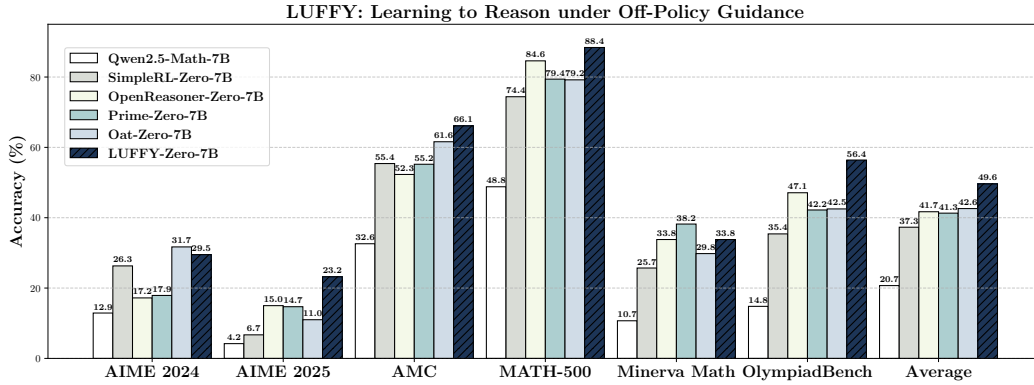


Figure 2: Overall performance across six competition-level benchmarks (AIME 2024, AIME 2025, AMC, MATH-500, Minerva Math, and OlympiadBench). **LUFFY** achieves an average score of **49.6**, delivering a substantial performance gain of over **+7.0** points on average compared to existing zero reinforcement learning methods.

1 Introduction

Recent breakthroughs in large reasoning models, including OpenAI-o1 [1], DeepSeek-R1 [2], and Kimi-1.5 [3], have demonstrated remarkable capabilities in complex reasoning tasks. These models have shown unprecedented proficiency in generating extensive Chains-of-Thought (CoT, [4]) responses and exhibiting sophisticated behaviors, such as self-reflection and self-correction. Particularly noteworthy is how these achievements have been realized through reinforcement learning with purely rule-based rewards, as demonstrated by DeepSeek-R1. The emergence of long CoT reasoning and self-reflection capabilities through such straightforward reward mechanisms, termed the “aha moment”, represents a significant advancement in the field.

One crucial paradigm behind the success is known as *zero-RL* [2, 5, 6, 7], which applies reinforcement learning to base language models directly, eliciting reasoning potentials using models’ own rollouts. Nevertheless, it has a fundamental limitation worth highlighting: it is inherently “on-policy”, constraining learning exclusively to the model’s self-generated outputs through iterative trials and feedback cycles. Despite showing promising results, zero-RL is bounded by the base LLM itself [8]. In essence, reinforcement learning under this setting amplifies existing behaviors rather than introducing genuinely novel cognitive capacities. Recent study [9] corroborates this constraint, demonstrating that models like Llama 3.2 [10] quickly reach performance plateaus under zero-RL training precisely because they lack certain foundational cognitive behaviors necessary for further advancement.

This inherent limitation provokes critical questions about the effectiveness and scope of learning within the zero-RL paradigm: *How can we empower LLMs to acquire reasoning behaviors surpassing their initial cognitive boundaries?* A natural approach to introduce external guidance from a stronger policy is imitation learning, wherein models are fine-tuned using reasoning traces produced by powerful LRMs like DeepSeek-R1 [2, 11, 12]. However, recent research [13, 14] raise concern about the generalization limits learned via pure imitation, which locks models into superficial and rigid reasoning models that impede further learning. Meanwhile, off-policy learning, which has proven powerful in various RL tasks for expanding learning beyond an agent’s initial capabilities [15, 16], remains largely unexplored in zero-RL. This leaves open critical questions about how to effectively incorporate off-policy knowledge alongside the exploration of on-policy learning, beyond simple imitation learning.

In this work, we aim to integrate off-policy guidance within the unified zero-RL paradigm, and introduce **LUFFY**: Learning to reason Under **off**-policy guidance. Based on conventional zero-RL methods such as GRPO [17], LUFFY introduces off-policy reasoning traces (e.g., from DeepSeek-R1) and combines them with models’ on-policy roll-outs before advantage computation, as illustrated in Figure 1. Intuitively, since off-policy traces consistently obtain positive rewards, LUFFY enables

the model to selectively imitate these high-quality reasoning traces when its own roll-outs fail to achieve correctness, while preserving the capacity for self-driven exploration whenever its generated reasoning steps are successful. In this way, LUFFY achieves a dynamic and adaptive equilibrium between imitation and exploration. However, naively combining off-policy traces can lead to overly rapid convergence and entropy collapse, causing the model to latch onto superficial patterns rather than acquiring genuine reasoning capabilities. To address these issues, we introduce *policy shaping via regularized importance sampling*, which amplifies learning signals for low-probability yet crucial actions under off-policy guidance. This mechanism encourages the model to preserve exploration throughout training, ultimately enabling it to internalize deeper and more generalizable reasoning behaviors.

As shown in Figure 2, LUFFY achieves clear improvements of **+7.0** points on average compared with previous RL-zero methods, across AIME24/25 [18], AMC [18], OlympiadBench [19], Minerva [20], and MATH-500 [21] benchmarks, establishing the effectiveness of off-policy learning in zero-like paradigms. Moreover, LUFFY demonstrates superior generalization capability, i.e., an advantage of over **+6.2** points on average, on out-of-distribution tasks, where SFT falls short. Empirical results demonstrate that LUFFY encourages the model to imitate high-quality reasoning traces while maintaining exploration of its own sampling space. This aligns with our deeper analysis, revealing that LUFFY assimilates off-policy reasoning behaviors flexibly and effectively, while SFT confines the model to rigid memorization of external reasoning patterns, hindering generalization and exploration.

2 Learning to Reason under Off-Policy Guidance

To facilitate exploration beyond the model’s own capabilities, we incorporate *off-policy guidance*, i.e., off-the-shelf reasoning trajectories generated by a stronger reasoning model such as Deepseek R1, into the zero-RL learning. We expect the model to learn generalizable knowledge from off-policy beyond superficial imitation and maintain effective and efficient exploration as in zero-RL training.

In the following sections, we first introduce the RL backbone algorithm, GRPO, followed by an illustration of mixed-policy GRPO, which naively integrates off-policy traces. Finally, we introduce LUFFY, leveraging policy shaping to mitigate entropy collapse and encourage continuous exploration.

2.1 GRPO and Importance Sampling

Due to the success of Deepseek-R1 [2], GRPO [17] becomes the de facto approach with zero-RL training. Compared to the widely used PPO [22], GRPO uses the reward scores of N sampled solutions from a query to estimate the advantage and thus remove the need for an additional value model. Formally, we denote the policy model before and after the update as $\pi_{\theta_{\text{old}}}$ and π_{θ} . Given a question q , a set of solutions τ_i generated by $\pi_{\theta_{\text{old}}}$, and the reward function $R(\cdot)$, the GRPO objective is defined as follows:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \frac{1}{\sum_{i=1}^N |\tau_i|} \sum_{i=1}^N \sum_{t=1}^{|\tau_i|} \min[r_{i,t}(\theta) A_i, \text{clip}(r_{i,t}(\theta); 1 - \epsilon, 1 + \epsilon) A_i] - \beta \cdot \mathbb{D}_{\text{KL}}[\pi_{\theta} \parallel \pi_{\text{ref}}]$$

$$\text{where } r_{i,t}(\theta) = \frac{\pi_{\theta}(\tau_{i,t} | q, \tau_{i,<t})}{\pi_{\theta_{\text{old}}}(\tau_{i,t} | q, \tau_{i,<t})}, A_i = \frac{R(\tau_i) - \text{mean}(\{R(\tau_i) \mid \tau_i \sim \pi_{\theta_{\text{old}}}(\tau), i = 1, 2, \dots, N\})}{\text{std}(\{R(\tau_i) \mid \tau_i \sim \pi_{\theta_{\text{old}}}(\tau), i = 1, 2, \dots, N\})} \quad (1)$$

\mathbb{D}_{KL} is the KL divergence. In the RL objective, GRPO follows PPO, using the importance sampling ($r_{i,t}$ in Eq.1) to calibrate the gradient as the rollouts are generated by $\pi_{\theta_{\text{old}}}$.

The clipping term with clip ratio ϵ empirically ensures that the current policy π_{θ} is within the trust region [23] of the old policy $\pi_{\theta_{\text{old}}}$. We loosely categorize this method as *On-Policy RL*, indicating that the model is optimized using samples drawn from distributions closely aligned with its current policy. Nevertheless, recent practices [24, 25, 7] have increasingly omitted the KL divergence term, making these methods somewhat “less On-Policy”.

2.2 Mixed-Policy GRPO

We incorporate off-policy rollouts in GRPO by adding them directly to the group of on-policy rollouts generated by the model itself. Provided an off-policy distribution π_ϕ , this would affect the advantage computations in the following way,

$$\hat{A}_i = \frac{R(\tau_i) - \text{mean}(\mathcal{G}_{\text{on}} \cup \mathcal{G}_{\text{off}})}{\text{std}(\mathcal{G}_{\text{on}} \cup \mathcal{G}_{\text{off}})}, \quad (2)$$

where $\mathcal{G}_{\text{on}} = \{R(\tau_i) \mid \tau_i \sim \pi_{\theta_{\text{old}}}(\tau), i = 1, 2, \dots, N_{\text{on}}\}$ and $\mathcal{G}_{\text{off}} = \{R(\tau_j) \mid \tau_j \sim \pi_\phi(\tau), j = 1, 2, \dots, N_{\text{off}}\}$. This group computation naturally assigns higher advantage to off-policy rollouts when the model struggles to generate correct solutions independently, while on-policy rollouts take precedence once the model begins producing successful reasoning traces, thereby encouraging self-driven exploration.

However, this mixed advantage computation introduces bias into the estimate of the policy gradient algorithm [26], which assumes the policy distribution generates the rollouts. In our preliminary experiments, this would lead to a huge performance drop in the training process. Therefore, we use importance sampling to calibrate gradient estimates and refer to this approach as *Mixed-Policy GRPO*:

$$\begin{aligned} \nabla_\theta \mathcal{J}(\theta) &= \mathbb{E}_{\tau_j \sim \pi_\theta(\tau)} \left[\nabla_\theta \log \pi_\theta(\tau_j) \hat{A}_j \right] \\ &= \mathbb{E}_{\tau_j \sim \pi_\phi(\tau)} \left[\frac{\pi_\theta(\tau_j)}{\pi_\phi(\tau_j)} \nabla_\theta \log \pi_\theta(\tau_j) \hat{A}_j \right]. \end{aligned} \quad (3)$$

The importance sampling term effectively corrects the gradient from expectation of π_θ (red) to π_ϕ (blue). This contrasts with the importance sampling term $r_{i,t}(\theta)$ used in on-policy RL (Eq. 1), where the denominator corresponds to the pre-update roll-out model policy $\pi_{\theta_{\text{old}}}$. Since the divergence between π_θ and $\pi_{\theta_{\text{old}}}$ is typically much smaller than that between π_θ and the off-policy policy π_ϕ , the importance sampling ratio in Eq. 3 tends to be smaller, serving to calibrate gradient estimates from a distinct distribution.

The RL objective of Mixed-Policy GRPO is extended from the original GRPO objective (Equation 1).

$$\begin{aligned} \mathcal{J}_{\text{Mixed}}(\theta) &= \frac{1}{Z} \underbrace{\left(\sum_{j=1}^{N_{\text{off}}} \sum_{t=1}^{|\tau_j|} \min[\hat{r}_{j,t}(\theta, \phi) \hat{A}_j, \text{clip}(\hat{r}_{j,t}(\theta, \phi); 1 - \epsilon, 1 + \epsilon) \hat{A}_j] \right)}_{\text{off-policy objective}} \\ &\quad + \underbrace{\sum_{i=1}^{N_{\text{on}}} \sum_{t=1}^{|\tau_i|} \min[r_{i,t}(\theta) \hat{A}_i, \text{clip}(r_{i,t}(\theta); 1 - \epsilon, 1 + \epsilon) \hat{A}_i]}_{\text{on-policy objective}}, \end{aligned} \quad (4)$$

$$\text{where } \hat{r}_{j,t}(\theta, \phi) = \frac{\pi_\theta(\tau_{j,t}|q, \tau_{j,<t})}{\pi_\phi(\tau_{j,t}|q, \tau_{j,<t})} \text{ and } r_{i,t}(\theta) = \frac{\pi_\theta(\tau_{i,t}|q, \tau_{i,<t})}{\pi_{\theta_{\text{old}}}(\tau_{i,t}|q, \tau_{i,<t})}.$$

$Z = \sum_{j=1}^{N_{\text{off}}} |\tau_j| + \sum_{i=1}^{N_{\text{on}}} |\tau_i|$ is the normalization factor. Based on the theoretical analysis of stochastic gradient descent in nonconvex optimization [27], we give a convergence analysis in Theorem 1 to show that our importance-weighted policy gradient estimator in Eq. (3) stabilizes and converges to a stationary point, and the convergence rate is $O(1/\sqrt{K})$, where K is the total number of iterations. The proof can be found in Appendix A.

Theorem 1. Suppose the objective function of the policy gradient algorithm $J \in \mathcal{J}_n$, where \mathcal{J}_n is the class of finite-sum Lipschitz smooth functions, has σ -bounded gradients, and the importance weight $w = \pi_\theta/\pi_\phi$ is clipped to be bounded by $[\underline{w}, \bar{w}]$. Let $\alpha_k = \alpha = c/\sqrt{K}$ where $c = \sqrt{\frac{2(J(\theta^*) - J(\theta^0))}{L\sigma^2 \underline{w}\bar{w}}}$, and θ^* is an optimal solution. Then, the iterates of our algorithm in Eq. (3) satisfy:

$$\min_{0 \leq k \leq K-1} \mathbb{E}[\|\nabla J(\theta^k)\|^2] \leq \sqrt{\frac{2(J(\theta^*) - J(\theta^0))L\bar{w}}{K\underline{w}}} \sigma.$$

The importance sampling ratio in off-policy learning typically involves π_ϕ , representing the behavior policy’s probability in off-policy trajectories [28]. Theoretically, our derivations and guarantees hold for any well-defined π_ϕ distribution. In practice, to facilitate direct integration of high-quality demonstrations from large, powerful models (e.g., DeepSeek-R1), we adopt $\pi_\phi = 1$ for computational efficiency. This practical choice avoids the complexity caused by different tokenization between the on-policy and off-policy models. It facilitates easy incorporation of off-the-shelf datasets without recomputation of π_ϕ , as well as preserving theoretical guarantees. In addition, we omit the clip operation for the off-policy rollouts, as the clip operation will be imbalanced when $\pi_\phi = 1$. The follow subsections illustrate **LUFFY**, which integrates policy shaping into Mixed-Policy GRPO.

2.3 Policy Shaping via Regularized Importance Sampling

While Mixed-Policy GRPO incorporates off-policy rollouts successfully via importance sampling, a new practical challenge emerges: importance sampling accelerates convergence but significantly reduces exploration (Figure 3, left). Specifically, entropy collapses much faster than in on-policy RL, indicating increasingly deterministic rollouts and a diminished capacity for exploring diverse reasoning trajectories.

This originates from the “hacking” of the Mixed-Policy objective. When combining both learning off-policy and on-policy signals, the model tends to quickly converge toward reinforcing off-policy tokens that are also likely in the on-policy π_θ distribution, and ignoring off-policy tokens that are deviated from the model’s original policy, i.e., low-probability tokens that may represent essential reasoning capabilities the model has yet to acquire. We empirically analyze this issue in detail in Section 4.2.

To address this issue, we introduce *policy shaping via regularized importance sampling*, a technique that re-weights the gradient of off-policy distributions to enhance learning from low-probability tokens. In particular, our approach replaces the importance sampling ratio $\pi_\theta(\tau_{j,t}|q, \tau_{j,<t})/\pi_\phi(\tau_{j,t}|q, \tau_{j,<t})$ with $f(\pi_\theta(\tau_{j,t}|q, \tau_{j,<t})/\pi_\phi(\tau_{j,t}|q, \tau_{j,<t}))$, where $f(\cdot)$ represents a transformation function that alters the dynamics between off-policy and on-policy distributions, thereby increasing gradient emphasis on tokens with low probability in the model’s standard distribution. Recall that we omit the clip operation for the off-policy rollouts. The loss function with policy shaping can be written as below:

$$\begin{aligned} \mathcal{J}_{\text{SHAPING}}(\theta) &= \frac{1}{Z} \left(\sum_{j=1}^{N_{\text{off}}} \sum_{t=1}^{|\tau_j|} f(\hat{r}_{j,t}(\theta, \phi)) \cdot \hat{A}_j \right. \\ &\quad \left. + \sum_{i=1}^{N_{\text{on}}} \sum_{t=1}^{|\tau_i|} \min[r_{i,t}(\theta) \cdot \hat{A}_i, \text{clip}(r_{i,t}(\theta); 1 - \epsilon, 1 + \epsilon) \cdot \hat{A}_i] \right), \end{aligned} \quad (5)$$

where $\hat{r}_{j,t}(\theta, \phi) = \frac{\pi_\theta(\tau_{j,t}|q, \tau_{j,<t})}{\pi_\phi(\tau_{j,t}|q, \tau_{j,<t})}$ and $r_{i,t}(\theta) = \frac{\pi_\theta(\tau_{i,t}|q, \tau_{i,<t})}{\pi_{\theta_{\text{old}}}(\tau_{i,t}|q, \tau_{i,<t})}$.

To further illustrate the meaning of shaping function f , we derive the gradient of off-policy objective as follows,

$$\begin{aligned} \nabla_\theta \mathcal{J}_{\text{SHAPING-OFF}}(\theta) &= \mathbb{E}_{\tau \sim \pi_\phi} \left[\nabla_\theta f\left(\frac{\pi_\theta}{\pi_\phi}\right) \cdot \hat{A}_j \right] \\ &= \mathbb{E}_{\tau \sim \pi_\phi} \left[f'\left(\frac{\pi_\theta}{\pi_\phi}\right) \frac{1}{\pi_\phi} \nabla_\theta \pi_\theta \cdot \hat{A}_j \right] \\ &= \mathbb{E}_{\tau \sim \pi_\phi} \left[\underbrace{f'(\pi_\theta) \frac{\pi_\theta}{\pi_\phi} \nabla_\theta \log \pi_\theta}_{\text{importance sampling}} \cdot \hat{A}_j \right]. \end{aligned} \quad (6)$$

We write $\pi(\tau_{j,t}|q, \tau_{j,<t})$ as π for simplicity. From the derivation, we can see $f'(\pi_\theta)$ is a weighting function of the gradient. The vanilla mixed-policy GRPO can be regarded as using a linear shaping function, i.e., $f(\pi) = \pi$, where the original importance sampling ratio π_θ/π_ϕ is applied.

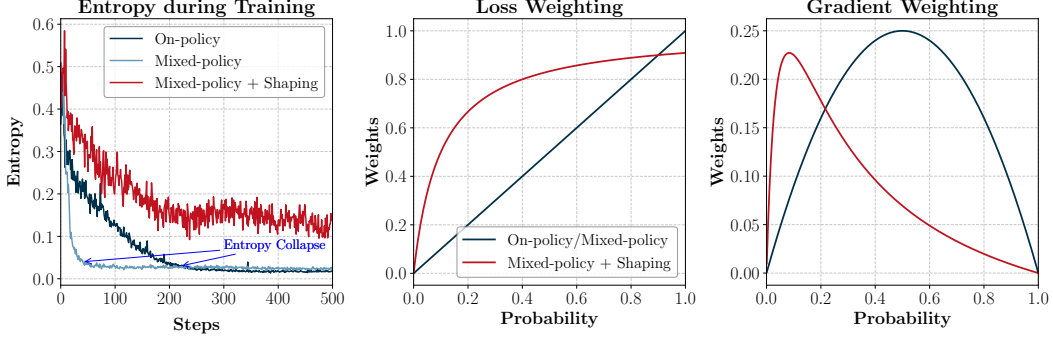


Figure 3: **Left:** generation entropy during training; **Middle:** shaping function value w.r.t. action probability; **Right:** gradient weights w.r.t. action probability.

For further insights on how $f'(\cdot)$ changes the model behavior, we can decompose the log probability and derive the gradient on each output logit:

$$\begin{aligned} \frac{\partial \mathcal{J}_{\text{SHAPING-OFF}}(\theta)}{\partial M_{\theta}(\tau'_{j,t})} &= \mathbb{E}_{\tau \sim \pi_{\phi}} \left[f'(\pi_{\theta}) \pi_{\theta} [\mathbb{I}(\tau'_{j,t} = \tau_{j,t}) - \pi_{\theta}] \cdot \hat{A}_j \right] \\ \Rightarrow \left| \frac{\partial \mathcal{J}_{\text{SHAPING-OFF}}(\theta)}{\partial M_{\theta}(\tau'_{j,t})} \right| &\leq \mathbb{E}_{\tau \sim \pi_{\phi}} \left[|f'(\pi_{\theta})| \pi_{\theta} (1 - \pi_{\theta}) \cdot |\hat{A}_j| \right], \end{aligned} \quad (7)$$

where $\tau'_{j,t}$ is any possible action/token on in the action space at the j -th trajectory and t -th position, and $M_{\theta}(\tau)$ denotes the logits of that action. The identity case represents the gradient when the action is the off-policy action $\tau = \tau'$, which elevates the probability of predicting the off-policy action. From Equation 7, we can see that the scale of gradient is upper-bounded by $\pi_{\theta}(1 - \pi_{\theta})$, leading to small values when $\pi_{\theta} \rightarrow 0$ and $\pi_{\theta} \rightarrow 1$. To encourage low-probability yet crucial actions, we use the $f(x) = x/(x + \gamma)$ as our shaping function (middle part of Figure 3), where γ is set as 0.1 for all experiments.

Considering the identity case ($\tau = \tau'$) and $\pi_{\phi} = 1$, the gradient with policy shaping can be written as:

$$\Rightarrow \mathbb{E}_{\tau \sim \pi_{\phi}} \left[\frac{\gamma}{(\pi_{\theta} + \gamma)^2} \pi_{\theta} (1 - \pi_{\theta}) \cdot \hat{A}_j \right]. \quad (8)$$

As shown in Figure 3 (Right), the shaping function reweights the gradients to assign more importance to low-probability actions, thereby improving learning from unfamiliar yet effective decisions from the off-policy traces.

2.4 Removing On-policy Clip

In PPO, the clipping mechanism is introduced to constrain policy updates within a trust region [23], thereby ensuring stable training. However, when incorporating off-policy guidance, the target behavior may deviate significantly from the model’s current policy, especially early in training.

As shown in Figure 4, LUFFY experiences more frequent clipping compared to On-Policy RL, which can suppress learning from high-quality off-policy traces. To address this, we remove the on-policy clip to allow greater flexibility in updating toward unfamiliar yet effective actions, thereby unlocking the model’s capacity to better integrate off-policy reasoning behaviors.

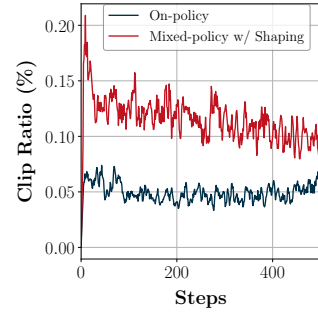


Figure 4: **Ratio of clipped signals.**

3 Experimental Setup

Dataset Construction. Our training set is a subset of OpenR1-Math-220k [29]³, of which the prompts are collected from NuminaMath 1.5 [30], and the off-policy reasoning traces are generated by Deepseek-R1 [2]. We use the default subset, which contains 94k prompts, and we filter out generations that are longer than 8192 tokens and those that are verified wrong by *Math-verify*⁴, resulting in 45k prompts and off-policy reasoning traces.

RL Practice. We remove the KL loss term by setting $\beta = 0$ and set entropy loss coefficient to 0.01 (Equation 1). Following Dr.GRPO[6], we remove the length normalization and standard error normalization of GRPO loss (Equation 1) for all experiments. For policy shaping, we empirically set the γ as 0.1. Our rollout batch size is 128, and the update batch size is 64. We use 8 samples per prompt. Specifically, for on-policy RL, we use 8 on-policy samples. For our methods, we use 1 off-policy and 7 on-policy samples to ensure fairness. We use temperature=1.0 for rollout generation.

Training. We use Qwen2.5-Math-7B [31] by default following previous work [24, 5, 6]. In addition to Qwen2.5-Math-7B, we extend LUFFY to Qwen2.5-Math-1.5B [31] and Qwen2.5-Instruct-7B [32]. Since the context length of Qwen2.5-Math is 4096 and the generation length of off-policy samples could be lengthy, we change the rope theta from 10000 to 40000 and extend the window size to 16384. The learning rate is constantly set as 1e-6. We train 500 steps for RL models. For the SFT baselines, we train all models for three epochs. All training experiments are conducted using 8 A100 GPUs. All our trained models share the same system prompt for training and inference:

System Prompt. Your task is to follow a systematic, thorough reasoning process before providing the final solution. This involves analyzing, summarizing, exploring, reassessing, and refining your thought process through multiple iterations. Structure your response into two sections: Thought and Solution. In the Thought section, present your reasoning using the format: “<think>\n thoughts </think>\n”. Each thought should include detailed analysis, brainstorming, verification, and refinement of ideas. After “</think>\n” in the Solution section, provide the final, logical, and accurate answer, clearly derived from the exploration in the Thought section. If applicable, include the answer in \boxed{ } for closed-form results like multiple choices or mathematical solutions.

Evaluation. For evaluation, we mainly focus on six widely used math reasoning benchmarks, including AIME 2024, AIME 2025, AMC [18], Minerva [20], OlympiadBench [19], and MATH-500 [21]. For AIME 2024, AIME 2025 and AMC, we report avg@32 as the test set is relatively small, and for the other three benchmarks, we report pass@1. As our RL training mainly focuses on math reasoning, we further validate the generalization capability on three out-of-distribution benchmarks, namely ARC-c [33](Open-Domain Reasoning), GPQA-diamond [34](Science Graduate Knowledge), and MMLU-Pro [35](Reasoning-focused Questions from Academic Exams and Textbooks). We shuffle the multiple-choice options to avoid contamination. For testing, the temperature is set as 0.6.

Reward Function. We use the rule-based reward, verified by Math-Verify.

$$R(q, \tau) = \begin{cases} 1 & \text{if } \tau \text{ outputs the correct final answer to } q \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

4 Experimental Results

Baseline Methods. For zero-RL methods, we consider the following methods: (1) *Simple-RL* [5]: training from Qwen2.5-Math-7B using rule-based reward; (2) *Oat-Zero* [6]: training from Qwen2.5-Math-7B and rule-based reward, proposing to remove the standard deviation in GRPO advantage computation and token-level normalization in policy loss computation; (3) *PRIME-Zero* [24]: using policy rollouts and outcome labels through implicit process rewards; (4) *OpenReasonerZero* [7]: a

³<https://huggingface.co/datasets/open-r1/OpenR1-Math-220k>

⁴<https://github.com/huggingface/Math-Verify>

Table 1: Overall performance on six competition-level reasoning benchmarks based on Qwen2.5-Math-7B-Base. Our method outperforms previous approaches with zero RL paradigm and achieves significant improvement over on-policy RL and SFT baselines. Bold and underline represent the 1st and 2nd in performance.

Model	AIME 24	AIME 25	AMC	MATH-500	Minerva	Olympiad	Avg.
Qwen2.5-Math [31]	12.9	4.2	32.6	48.8	10.7	14.8	20.7
Qwen2.5-Math-Instruct [31]	11.4	8.8	48.3	81.2	33.1	38.8	36.9
SimpleRL-Zero [5]	26.3	6.7	55.4	74.4	25.7	35.4	37.3
OpenReasoner-Zero [7]	17.2	15.0	52.3	84.6	33.8	47.1	41.7
PRIME-Zero [24]	17.9	14.7	55.2	79.4	38.2	42.2	41.3
Oat-Zero [6]	31.7	11.0	<u>61.6</u>	79.2	29.8	42.5	42.6
SFT	28.6	23.5	59.0	86.0	37.5	51.1	47.6
On-Policy RL	24.6	15.7	61.3	84.6	34.9	47.9	44.8
LUFFY	<u>29.5</u>	<u>23.2</u>	66.1	88.4	33.8	56.4	49.6

Table 2: Overall out-of-distribution performance on three benchmark datasets (based on Qwen2.5-Math-7B-Base).

Model	ARC-c	GPQA-diamond	MMLU-Pro	Avg.
Qwen2.5-Math-7B-Base [31]	18.2	11.1	16.9	15.4
Qwen2.5-Math-7B-Instruct [31]	70.3	24.7	34.1	43.0
SimpleRL-Zero [5]	30.2	23.2	34.5	29.3
OpenReasoner-Zero [7]	66.2	29.8	58.7	51.6
PRIME-Zero [24]	73.3	18.2	32.7	41.4
Oat-Zero [6]	70.1	23.7	41.7	45.2
SFT	75.2	24.7	42.7	47.5
On-Policy RL	82.3	40.4	49.3	<u>57.3</u>
LUFFY	<u>80.5</u>	<u>39.9</u>	<u>53.0</u>	57.8

recent open-source implementation of zero-RL methods. Except zero-RL approaches from previous work, we consider two more baselines with our setting (1) *On-Policy RL* – we train on-policy RL with zero-RL paradigm using Dr.GRPO with the same reward and data. (2) *SFT* – we train the model with the same prompts and reasoning traces as LUFFY using SFT.

4.1 Main Results

Reasoning Benchmark Performance. Our main results are presented in Table 1, comparing LUFFY against representative zero-RL baselines of similar computational complexity, including our reproduced On-Policy RL. All prior methods are built upon Qwen-7B base models, differing in dataset composition (source and difficulty) and optimization strategies, e.g., removing length and standard error normalization [6] or incorporating process-level rewards [24]. Evaluated on six challenging competition-level benchmarks, LUFFY achieves an average score of **49.6**, outperforming existing zero-RL methods by a substantial margin of **+7.0** points, establishing a new state-of-the-art. Notably, while LUFFY exhibits comparable performance in AIME 24, it demonstrates a significantly greater advantage on the newly released AIME 25 test set (+8.2), demonstrating its generalization to internalize nuanced reasoning behaviors from off-policy traces. Compared to our zero-RL replication, i.e., On-Policy RL, LUFFY improves performance by +4.8 points on average, clearly demonstrating the benefit of integrating high-quality off-policy traces into the RL training loop. Moreover, compared to its SFT counterpart, our method achieves a notable improvement of +2.0 points on average, providing a more robust and effective alternative for *distilling knowledge* from stronger LRMs, traditionally achieved through supervised fine-tuning [2, 11, 12].

Out-of-Distribution Generalization. We further investigate the generalization capabilities of our method beyond mathematical reasoning, using three challenging out-of-distribution benchmarks, as summarized in Table 2. LUFFY demonstrates considerable performance gains over previous zero-RL practices. Although SFT achieves competitive results on mathematical reasoning tasks (Table 1), it struggles to generalize effectively to domains significantly different from its training distribution,

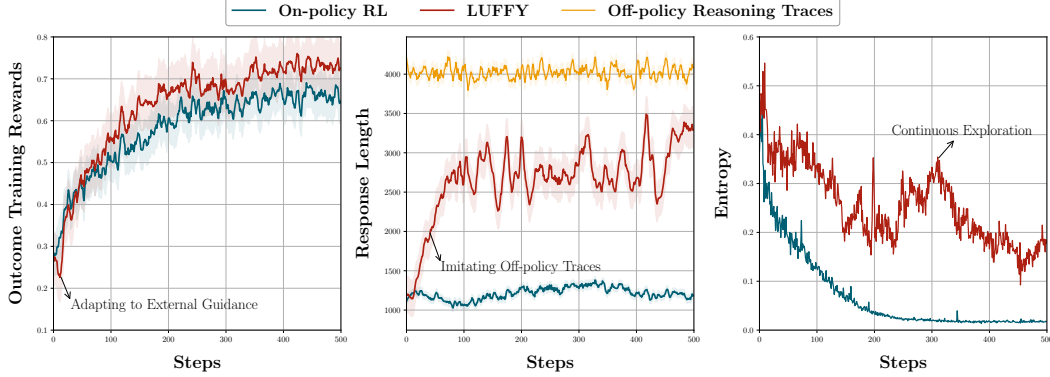


Figure 5: Training dynamics of LUFFY compared with on-policy RL. **Left:** outcome training rewards; **Middle:** generation length; **Right:** generation entropy.

such as open-domain reasoning (ARC-C), graduate-level science knowledge (GPQA-diamond), and general academic knowledge (MMLU-Pro). In contrast, LUFFY matches the strong performance of On-Policy RL on ARC-C and GPQA-diamond and substantially outperforms it on the MMLU-Pro benchmark, achieving comparable average results. Notably, compared to Qwen2.5-Math-7B-Instruct, LUFFY achieves considerable improvements in both specialized reasoning tasks (+12.7 points on average, Table 1) and OOD tasks (+14.8 points on average), all while utilizing significantly less training data and computational resources. These findings underscore the effectiveness of LUFFY in leveraging off-policy reasoning guidance for enhanced generalization across diverse, out-of-distribution tasks.

4.2 Training Dynamics of LUFFY

Strategically Learning from Guidance. Figure 5 illustrate the training dynamics regarding training rewards, generation length and entropy for On-Policy RL and LUFFY. Initially, LUFFY primarily imitates off-policy trajectories, as indicated by the increasing generation length gradually aligning with the off-policy reasoning traces (middle part of Figure 5). At this early stage, imitation dominates, causing an initial performance dip (left part of Figure 5) as the model adjusts to external guidance. As training progresses, on-policy rollouts gradually become more prominent, fostering independent exploration within the model’s own sampling space while effectively retaining insights gained from off-policy demonstrations. This guided exploration brings growing advantages (training rewards) over On-Policy RL. Consequently, LUFFY achieves a dynamic balance between imitation and exploration, leading to more effective off-policy learning (Section 5). These results highlight that LUFFY selectively adopts valuable reasoning patterns rather than blindly imitating off-policy traces.

Maintaining Exploration. Figure 5 (Right) illustrates that LUFFY consistently sustains higher entropy compared to On-Policy RL throughout the entire training process. Specifically, the generation entropy of On-Policy RL rapidly converges to nearly zero after approximately 200 steps, indicating a highly deterministic policy with limited exploration potential. Conversely, the elevated entropy observed in LUFFY allows continuous exploration of less confident yet potentially superior policies, facilitating the discovery and learning of novel cognitive behaviors. Interestingly, we observe entropy fluctuations and even occasional increases, such as between steps 200 and 250, reflecting ongoing exploration of low-probability but crucial actions, also referred to as *pivotal tokens* [36, 25]. This strategic exploration enables the model to escape local optima, thus improving its convergence towards more globally optimal solutions.

4.3 Ablation Study

In this section, we perform an ablation study to examine the contributions of LUFFY components, as summarized in Table 3. Shaping and NoClip both positively contribute to the final performance of Mixed-Policy training. However, applying these enhancements without

Table 3: Ablation study on LUFFY components.

Model	AIME 24	AIME 25	AMC	MATH-500	Minerva	Olympiad	Avg.
Mixed-Policy RL	19.2	16.9	58.7	83.8	30.9	49.9	43.2
+ Shaping	30.0	22.6	61.7	86.2	36.4	55.6	48.7
+ Shaping + NoClip	29.5	23.2	66.1	88.4	33.8	56.4	49.6
On-Policy RL	24.6	15.7	61.3	84.6	34.9	47.9	44.8
+ Shaping	21.8	14.7	57.6	81.6	33.5	44.6	42.3
+ No Clip	22.7	17.3	61.6	83.4	34.9	50.7	45.1

Table 4: Overall performance on six competition-level benchmark performance on Qwen2.5-Math-1.5B and Qwen2.5-Instruct-7B.

Model	AIME 24	AIME 25	AMC	MATH-500	Minerva	Olympiad	Avg.
<i>Qwen2.5-Math-1.5B</i>							
Qwen2.5-Math-1.5B-Base [31]	7.9	4.7	26.4	31.0	12.1	21.5	17.3
Qwen2.5-Math-1.5B-Instruct [31]	11.4	8.5	47.4	75.2	27.6	38.7	34.8
SFT	15.2	14.3	43.5	74.8	30.9	36.9	40.3
On-Policy RL	12.6	6.5	42.6	68.8	22.1	34.4	36.1
LUFFY	15.2	12.7	46.8	79.4	26.5	42.4	42.1
<i>Qwen2.5-Instruct-7B</i>							
Qwen2.5-7B-Instruct [32]	11.9	7.6	44.1	74.6	30.5	39.7	34.7
SFT	9.7	11.3	41.8	71.2	26.8	38.1	33.1
On-Policy RL	16.5	10.7	47.5	75.8	35.3	41.9	37.9
LUFFY	16.6	15.7	52.2	81.4	36.8	48.7	41.9

off-policy guidance (On-Policy + No Clip/Shaping) does not yield improvement, underscoring the necessity of external signals to acquire nuanced and generalizable reasoning skills. Additionally, Figure 6 illustrates the validation performance throughout training to provide insights into the training dynamics. Notably, Mixed-Policy demonstrates rapid initial improvement, substantially outperforming On-Policy RL during early stages. However, as training progresses, its performance converges to that of On-Policy RL.

These observations align with our earlier discussion on entropy collapse (Section 2.3), highlighting that directly integrating off-policy traces accelerates convergence but fails to prevent the model from becoming trapped in local optima. Conversely, policy shaping acts as an effective regularizer, mitigating premature convergence and consistently enlarging performance advantages in later training phases. This benefit is further amplified when removing the on-policy clip, enabling parameter updates that support more aggressive and effective exploration.

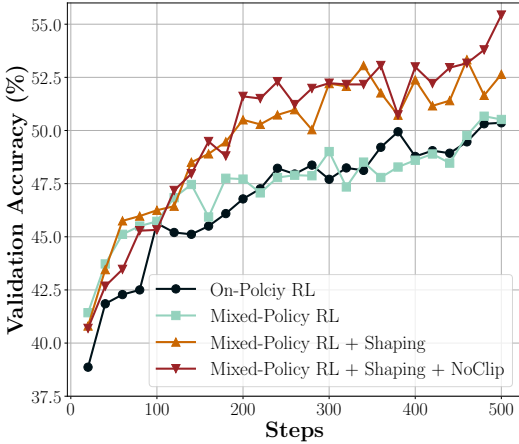


Figure 6: Validation performance during training.

4.4 Extension to More Models

To assess the general applicability of our method, we extend LUFFY to a smaller base model: Qwen2.5-Math-1.5B. We also consider Qwen2.5-Instruct-7B to validate LUFFY on non-zero RL paradigm, i.e., RL on SFT models. Table 4 presents the performance across six challenging competition-level benchmarks. On both models, LUFFY achieves consistent and substantial improvements, surpassing both SFT and On-Policy RL. On Qwen2.5-Math-1.5B, LUFFY attains an average score of **42.1**, demonstrating notable gains of +1.8 and +6.0 points over SFT and On-Policy RL, respectively.

Similar advantages are observed on the Qwen2.5-Instruct-7B model, where LUFFY consistently outperforms baselines across all benchmarks. Particularly, it achieves a overall score of **41.9**, representing clear improvements of +8.8 points over SFT and +4.0 points compared with On-Policy RL.

5 Analysis

In this section, we analyze how LUFFY effectively leverages off-policy guidance, i.e., *imitating to illuminate*, to improve reasoning quality and generalization.

LUFFY Selectively Learns Off-policy Traces. Imitation behavior can be observed through the similarity between model outputs and off-policy traces. To quantify this, we compare generations from SFT, On-Policy RL, and LUFFY against those from DeepSeek-R1 on a held-out set of 1,000 samples, using BLEU [37] as the similarity metric. The resulting BLEU scores are 57.5 for SFT, 8.8 for On-Policy RL, and 44.8 for LUFFY, reflecting the strong imitation behavior of SFT and the more selective, yet substantial, imitation in LUFFY.

LUFFY Learns Off-policy Traces More Effectively. We compare the generation length distributions of LUFFY and SFT on the combined set of six mathematical reasoning benchmarks. As shown in Figure 7, LUFFY produces significantly shorter generations on average (2,832 tokens) compared to SFT (4,646 tokens), suggesting a more effective reasoning process that balances imitation and exploration. In contrast, SFT often mimics the surface form of off-policy demonstrations without genuinely engaging in problem-solving. This behavior is especially evident in incorrect outputs, where SFT frequently generates overly long and ultimately unproductive reasoning traces. These results indicate that while both methods are exposed to similar off-policy signals, LUFFY learns to selectively internalize useful reasoning patterns, whereas SFT tends to overfit to superficial features of the off-policy data. We present a case study in Appendix B.

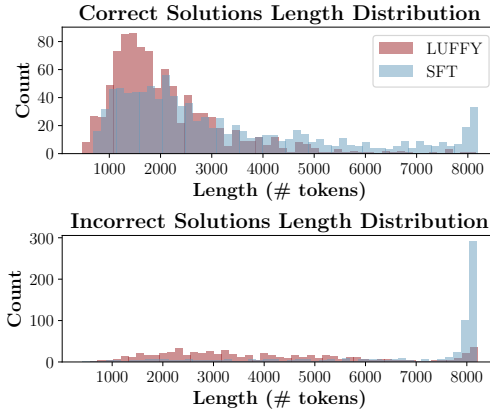


Figure 7: Generation length of correct and incorrect solutions.

LUFFY Can Explore During Test-time While SFT Cannot. We compute pass@8 accuracy on the combined AIME 2024 and AMC datasets, varying the generation temperature from 0.1 to 1.0. As shown in Figure 8, both RL-based methods (On-Policy RL and LUFFY) exhibit strong exploratory capabilities, with pass@8 improving as the temperature increases, showing potentials in scaling test-time compute [38]. In contrast, although SFT performs comparably to LUFFY under near-deterministic decoding (temperature 0.1), its performance deteriorates at higher temperatures, failing to uncover additional correct reasoning paths. This highlights the fragility and limited adaptability of SFT, which aligns with prior findings [13, 14] that suggest SFT tends to memorize reasoning patterns rather than learning generalizable reasoning capability.

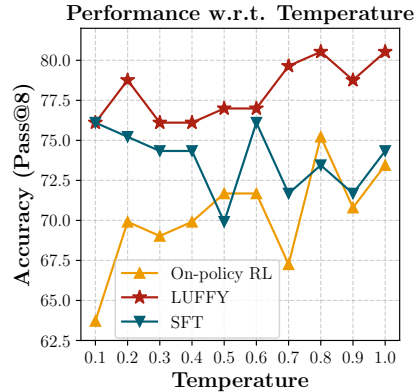


Figure 8: Pass@8 accuracy (on the merge sets of AIME 2024 and AMC) under different generation temperatures.

6 Related Work

RL for LRMs Recent advances have demonstrated remarkable progress in enhancing LLMs’ reasoning capabilities through reinforcement learning approaches [2, 1, 3, 39]. DeepSeek R1 pioneered a simple yet effective rule-based reward model, while OpenAI’s O1 advanced the field with step-by-step natural language feedback for precise policy optimization. Kimi 1.5 further contributed by introducing long-to-short optimization techniques for improved reasoning efficiency. Subsequent work has expanded the frontiers of chain-of-thought reasoning. BOLT [40] presented a novel three-stage framework for bootstrapping long chain-of-thought capabilities without relying on model distillation. LIMO [12] challenged conventional wisdom by showing that complex reasoning emerges from minimal, strategically curated examples. Concise Reasoning [41] further advanced the field through innovative RL-guided techniques for reasoning optimization. Significant methodological contributions have emerged in training and optimization approaches. TRPA [42] introduced a robust RL algorithm specifically designed for reasoning enhancement, while SimpleRL-Zoo [43] provided a comprehensive empirical study of zero-shot RL applications. Light-r1 [44] demonstrated the effectiveness of combining curriculum learning with DPO, and ThinkPO [45] introduced novel mechanisms for reasoning preference alignment. However, these methods either use supervised fine-tuning which shows limited generalization compared to RL approaches, or employ pure RL approaches that face efficiency challenges in exploration. Our work leverages high-quality off-policy data to enhance RL exploration efficiency while preserving the benefits of on-policy learning.

On-Policy and Off-Policy RL Reinforcement learning methods can be broadly categorized into on-policy and off-policy approaches based on how they use collected experiences for policy updates. On-policy methods, including TRPO [23], A2C/A3C [46], ACKTR [47], and PPO [22], update policies using only data collected from the current policy, ensuring stability but potentially limiting sample efficiency. Off-policy methods such as DQN [48], DDPG [49], TD3 [50], and SAC [51] can learn from experiences collected by different policies, enabling better sample efficiency but often at the cost of increased complexity and potential instability. In LLM training, on-policy methods are more commonly adopted, with approaches like GRPO [17], REINFORCE [52], ReMax [53], and PPO [22] demonstrating strong performance through various optimization techniques. Meanwhile, some works explore off-policy learning, such as DPO [54] which reformulates preference optimization as a classification problem. However, few work have investigated how to effectively combine the stability benefits of on-policy learning with the capability expansion potential of off-policy data. Our work bridges this gap by proposing novel techniques to leverage existing high-quality off-policy data to enhance on-policy learning while maintaining training stability.

7 Conclusion

We presented **LUFFY**, a simple yet powerful framework that integrates off-policy reasoning guidance into the zero-RL paradigm. By dynamically balancing imitation and exploration, LUFFY effectively leverages external reasoning traces without sacrificing the model’s ability to discover novel solutions. Our method outperforms strong baselines across competitive math benchmarks and generalizes robustly to out-of-distribution tasks, surpassing both on-policy RL and supervised fine-tuning. These results highlight the promise of off-policy learning as a scalable and principled path toward building more general, capable, and self-improving reasoning models. Future work may focus on extending LUFFY to broader domains or modalities [55] and further refining policy shaping to maximize exploration under off-policy guidance.

References

- [1] 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.
- [2] 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.
- [3] 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.
- [4] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [5] Weihao Zeng, Yuzhen Huang, Wei Liu, Keqing He, Qian Liu, Zejun Ma, and Junxian He. 7b model and 8k examples: Emerging reasoning with reinforcement learning is both effective and efficient. <https://hkust-nlp.notion.site/simpler1-reason>, 2025. Notion Blog.
- [6] 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*, 2025.
- [7] Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xiangyu Zhang, and Heung-Yeung Shum. Open-reasoner-zero: An open source approach to scaling up reinforcement learning on the base model, 2025.
- [8] Rosie Zhao, Alexandru Meterez, Sham Kakade, Cengiz Pehlevan, Samy Jelassi, and Eran Malach. Echo chamber: RL post-training amplifies behaviors learned in pretraining, 2025.
- [9] 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, 2025.
- [10] Meta AI. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models, September 2024. 15 minute read.
- [11] Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025.
- [12] Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie Xia, and Pengfei Liu. Limo: Less is more for reasoning. *arXiv preprint arXiv:2502.03387*, 2025.
- [13] Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V. Le, Sergey Levine, and Yi Ma. Sft memorizes, rl generalizes: A comparative study of foundation model post-training, 2025.
- [14] Hardy Chen, Haoqin Tu, Fali Wang, Hui Liu, Xianfeng Tang, Xinya Du, Yuyin Zhou, and Cihang Xie. Sft or rl? an early investigation into training rl-like reasoning large vision-language models. <https://github.com/UCSC-VLAA/VLAA-Thinking>, 2025.
- [15] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [16] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- [17] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024.

- [18] Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif Rasul, Longhui Yu, Albert Q. Jiang, Ziju Shen, et al. Numinamath: The largest public dataset in ai4maths with 860k pairs of competition math problems and solutions. <https://huggingface.co/datasets/Numinamath>, 2024. Hugging Face repository, 13:9.
- [19] Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3828–3850, 2024.
- [20] Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. *Advances in Neural Information Processing Systems*, 35:3843–3857, 2022.
- [21] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.
- [22] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [23] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR, 2015.
- [24] Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang, Wendi Li, Bingxiang He, Yuchen Fan, Tianyu Yu, Qixin Xu, Weize Chen, et al. Process reinforcement through implicit rewards. *arXiv preprint arXiv:2502.01456*, 2025.
- [25] Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, Weinan Dai, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. Dapo: An open-source llm reinforcement learning system at scale, 2025.
- [26] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- [27] Sashank J Reddi, Ahmed Hefny, Suvrit Sra, Barnabas Póczos, and Alex Smola. Stochastic variance reduction for nonconvex optimization. In *ICML*, pages 314–323, 2016.
- [28] Wenjia Meng, Qian Zheng, Gang Pan, and Yilong Yin. Off-policy proximal policy optimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 9162–9170, 2023.
- [29] Hugging Face. Open r1: A fully open reproduction of deepseek-r1, January 2025.
- [30] 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. [<https://huggingface.co/AI-MO/NuminaMath-1.5>] (https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf), 2024.
- [31] An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement, 2024.
- [32] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan,

- Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [33] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*, 2018.
 - [34] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.
 - [35] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhnanil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*, 2024.
 - [36] Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J. Hewett, Mojan Javaheripi, Piero Kauffmann, James R. Lee, Yin Tat Lee, Yuanzhi Li, Weishung Liu, Caio C. T. Mendes, Anh Nguyen, Eric Price, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Xin Wang, Rachel Ward, Yue Wu, Dingli Yu, Cyril Zhang, and Yi Zhang. Phi-4 technical report, 2024.
 - [37] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In Pierre Isabelle, Eugene Charniak, and Dekang Lin, editors, *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguistics.
 - [38] Amrith Setlur, Nived Rajaraman, Sergey Levine, and Aviral Kumar. Scaling test-time compute without verification or RL is suboptimal. In *ICLR 2025 Workshop: VerifAI: AI Verification in the Wild*, 2025.
 - [39] Xiaoye Qu, Yafu Li, Zhaochen Su, Weigao Sun, Jianhao Yan, Dongrui Liu, Ganqu Cui, Daizong Liu, Shuxian Liang, Junxian He, Peng Li, Wei Wei, Jing Shao, Chaochao Lu, Yue Zhang, Xian-Sheng Hua, Bowen Zhou, and Yu Cheng. A survey of efficient reasoning for large reasoning models: Language, multimodality, and beyond, 2025.
 - [40] Bo Pang, Hanze Dong, Jiacheng Xu, Silvio Savarese, Yingbo Zhou, and Caiming Xiong. Bolt: Bootstrap long chain-of-thought in language models without distillation. *arXiv preprint arXiv:2502.03860*, 2025.
 - [41] Mehdi Fatemi, Banafsheh Rafiee, Mingjie Tang, and Kartik Talamadupula. Concise reasoning via reinforcement learning. *arXiv preprint arXiv:2504.05185*, 2025.
 - [42] Xuerui Su, Shufang Xie, Guoqing Liu, Yingce Xia, Renqian Luo, Peiran Jin, Zhiming Ma, Yue Wang, Zun Wang, and Yuting Liu. Trust region preference approximation: A simple and stable reinforcement learning algorithm for llm reasoning. *arXiv preprint arXiv:2504.04524*, 2025.
 - [43] Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the wild. *arXiv preprint arXiv:2503.18892*, 2025.
 - [44] Liang Wen, Yunke Cai, Fenrui Xiao, Xin He, Qi An, Zhenyu Duan, Yimin Du, Junchen Liu, Lifu Tang, Xiaowei Lv, et al. Light-rl: Curriculum sft, dpo and rl for long cot from scratch and beyond. *arXiv preprint arXiv:2503.10460*, 2025.
 - [45] Wang Yang, Hongye Jin, Jingfeng Yang, Vipin Chaudhary, and Xiaotian Han. Thinking preference optimization. *arXiv preprint arXiv:2502.13173*, 2025.
 - [46] 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*, pages 1928–1937. PmLR, 2016.
 - [47] Yuhuai Wu, Elman Mansimov, Roger B Grosse, Shun Liao, and Jimmy Ba. Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation. *Advances in neural information processing systems*, 30, 2017.

- [48] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [49] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [50] Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pages 1587–1596. PMLR, 2018.
- [51] 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*, pages 1861–1870. Pmlr, 2018.
- [52] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- [53] Ziniu Li, Tian Xu, Yushun Zhang, Zhihang Lin, Yang Yu, Ruoyu Sun, and Zhi-Quan Luo. Remax: A simple, effective, and efficient reinforcement learning method for aligning large language models. *arXiv preprint arXiv:2310.10505*, 2023.
- [54] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.
- [55] Yunzhuo Hao, Jiawei Gu, Huichen Will Wang, Linjie Li, Zhengyuan Yang, Lijuan Wang, and Yu Cheng. Can mllms reason in multimodality? emma: An enhanced multimodal reasoning benchmark. *arXiv preprint arXiv:2501.05444*, 2025.
- [56] Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2 edition, 2018.

A Convergence Rate of the Importance-Weighted Policy Gradient Estimator

We study the nonconvex *finite-sum* problems of the form

$$\max_{\boldsymbol{\theta} \in \mathbb{R}^d} J(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n J_i(\boldsymbol{\theta}), \quad (10)$$

where both J and J_i ($i \in [n]$) may be nonconvex. We denote the class of such finite-sum Lipschitz smooth functions by $J \in \mathcal{J}_n$. Here, we optimize functions in \mathcal{J}_n of our importance-weighted policy gradient estimator.

The vanilla policy gradient algorithm maximizes the expected advantage function (equivalent to minimizing the negative expected advantage function) as

$$\max_{\boldsymbol{\theta} \in \mathbb{R}^d} J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim \pi_{\boldsymbol{\theta}}} [A(\tau)] \approx \frac{1}{n} \sum_{i=1}^n [A(\tau_i)], \quad (11)$$

According to the Policy Gradient Theorem [56], the vanilla policy gradient estimator has the following form:

$$\nabla J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim \pi_{\boldsymbol{\theta}}} [\nabla \log \pi_{\boldsymbol{\theta}}(\tau) \cdot A(\tau)] \approx \frac{1}{n} \sum_{i=1}^n [\nabla \log \pi_{\boldsymbol{\theta}}(\tau_i) \cdot A(\tau_i)], \quad (12)$$

where we use $\nabla J(\boldsymbol{\theta})$ to denote $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$ for simplicity. Our algorithm draws samples from another behavior policy π_{ϕ} , resulting in an importance-weighted policy gradient estimator as

$$\tilde{\nabla} J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim \pi_{\phi}} \left[\frac{\pi_{\boldsymbol{\theta}}(\tau_i)}{\pi_{\phi}(\tau_i)} \cdot \nabla \log \pi_{\boldsymbol{\theta}}(\tau) \cdot A(\tau) \right] \approx \frac{1}{n} \sum_{i=1}^n [w_i \cdot \nabla J_i(\boldsymbol{\theta})], \quad (13)$$

where $w_i = \frac{\pi_{\boldsymbol{\theta}}(\tau_i)}{\pi_{\phi}(\tau_i)}$ is the importance weight assigned to sample i .

Let α_k denote the learning rate at iteration k , and w_{i_k} be the instance weight assigned to sample i by our algorithm. By stochastic gradient ascent, our algorithm has the following update rule:

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k + \alpha_k w_{i_k} \nabla J_{i_k}(\boldsymbol{\theta}^k), i \in [n]. \quad (14)$$

Definition 1. For $J \in \mathcal{J}_n$, our algorithm takes an index $i \in [n]$ and a point $x \in \mathbb{R}^d$, and returns the pair $(J_i(\boldsymbol{\theta}), \nabla J_i(\boldsymbol{\theta}))$.

Definition 2. We say $J : \mathbb{R}^d \rightarrow \mathbb{R}$ is Lipschitz smooth (L -smooth) if there is a constant L such that

$$\|\nabla J(\boldsymbol{\vartheta}) - \nabla J(\boldsymbol{\theta})\| \leq L \|\boldsymbol{\vartheta} - \boldsymbol{\theta}\|, \quad \forall \boldsymbol{\vartheta}, \boldsymbol{\theta} \in \mathbb{R}^d. \quad (15)$$

Definition 3. A point $\boldsymbol{\theta}$ is called ϵ -accurate if $\|\nabla J(\boldsymbol{\theta})\|^2 \leq \epsilon$. A stochastic iterative algorithm is said to achieve ϵ -accuracy in k iterations if $\mathbb{E}[\|\nabla J(\boldsymbol{\theta}^k)\|^2] \leq \epsilon$, where the expectation is over the stochasticity of the algorithm.

Definition 4. We say $J \in \mathcal{J}_n$ has σ -bounded gradients if $\|\nabla J_i(\boldsymbol{\theta})\| \leq \sigma$ for all $i \in [n]$ and $\boldsymbol{\theta} \in \mathbb{R}^d$.

Definition 5. We say the positive instance weight w in our algorithm is bounded if there exist constants \underline{w} and \overline{w} such that $\underline{w} \leq w_i \leq \overline{w}$ for all $i \in [n]$.

Theorem 1. Suppose the objective function of the policy gradient algorithm $J \in \mathcal{J}_n$, where \mathcal{J}_n is the class of finite-sum Lipschitz smooth functions, has σ -bounded gradients, and the importance weight $w = \pi_{\boldsymbol{\theta}}/\pi_{\phi}$ is clipped to be bounded by $[\underline{w}, \overline{w}]$. Let $\alpha_k = \alpha = c/\sqrt{K}$ where $c = \sqrt{\frac{2(J(\boldsymbol{\theta}^*) - J(\boldsymbol{\theta}^0))}{L\sigma^2 \underline{w} \overline{w}}}$, and $\boldsymbol{\theta}^*$ is an optimal solution. Then, the iterates of our algorithm in Eq. (3) satisfy:

$$\min_{0 \leq k \leq K-1} \mathbb{E}[\|\nabla J(\boldsymbol{\theta}^k)\|^2] \leq \sqrt{\frac{2(J(\boldsymbol{\theta}^*) - J(\boldsymbol{\theta}^0))L\overline{w}}{K\underline{w}}} \sigma.$$

Proof. According to the Lipschitz continuity of ∇J , the iterates of our algorithm satisfy the following bound:

$$\mathbb{E}[J(\boldsymbol{\theta}^{k+1})] \geq \mathbb{E}[J(\boldsymbol{\theta}^k) + \langle \nabla J(\boldsymbol{\theta}^k), \boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^k \rangle - \frac{L}{2} \|\boldsymbol{\theta}^{k+1} - \boldsymbol{\theta}^k\|^2]. \quad (16)$$

After substituting (14) into (16), we have:

$$\begin{aligned}\mathbb{E}[J(\boldsymbol{\theta}^{k+1})] &\geq \mathbb{E}[J(\boldsymbol{\theta}^k)] + \alpha_k w_k \mathbb{E}[\|\nabla J(\boldsymbol{\theta}^k)\|^2] - \frac{L\alpha_k^2 w_k^2}{2} \mathbb{E}[\|\nabla J_{i_k}(\boldsymbol{\theta}^k)\|^2] \\ &\geq \mathbb{E}[J(\boldsymbol{\theta}^k)] + \alpha_k w_k \mathbb{E}[\|\nabla J(\boldsymbol{\theta}^k)\|^2] - \frac{L\alpha_k^2 w_k^2}{2} \sigma^2.\end{aligned}\quad (17)$$

The first inequality follows from the unbiasedness of the stochastic gradient $\mathbb{E}_{i_t}[\nabla J_{i_k}(\boldsymbol{\theta}^k)] = \nabla J(\boldsymbol{\theta}^k)$. The second inequality uses the assumption on gradient boundedness in Definition 4. Re-arranging (17) we obtain

$$\mathbb{E}[\|\nabla J(\boldsymbol{\theta}^k)\|^2] \leq \frac{1}{\alpha_k w_k} \mathbb{E}[J(\boldsymbol{\theta}^{k+1}) - J(\boldsymbol{\theta}^k)] + \frac{L\alpha_k w_k}{2} \sigma^2. \quad (18)$$

Summing (18) from $k = 0$ to $K - 1$ and using that α_k is fixed α we obtain

$$\begin{aligned}\min_t \mathbb{E}[\|\nabla J(\boldsymbol{\theta}^k)\|^2] &\leq \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla J(\boldsymbol{\theta}^k)\|^2] \\ &\leq \frac{1}{K} \sum_{k=0}^{K-1} \frac{1}{\alpha w_k} \mathbb{E}[J(\boldsymbol{\theta}^{k+1}) - J(\boldsymbol{\theta}^k)] + \frac{1}{K} \sum_{k=0}^{K-1} \frac{L\alpha w_k}{2} \sigma^2 \\ &\leq \frac{1}{K\alpha\bar{w}} (J(\boldsymbol{\theta}^K) - J(\boldsymbol{\theta}^0)) + \frac{L\alpha\bar{w}}{2} \sigma^2 \\ &\leq \frac{1}{K\alpha\bar{w}} (J(\boldsymbol{\theta}^*) - J(\boldsymbol{\theta}^0)) + \frac{L\alpha\bar{w}}{2} \sigma^2 \\ &\leq \frac{1}{\sqrt{K}} \left(\frac{1}{c\bar{w}} (J(\boldsymbol{\theta}^*) - J(\boldsymbol{\theta}^0)) + \frac{Lc\bar{w}}{2} \sigma^2 \right).\end{aligned}\quad (19)$$

The first step holds because the minimum is less than the average. The second step is obtained from (18). The third step follows from the assumption on instance weight boundedness in Definition 5. The fourth step is obtained from the fact that $J(\boldsymbol{\theta}^*) \geq J(\boldsymbol{\theta}^K)$. The final inequality follows upon using $\alpha = c/\sqrt{K}$. By setting

$$c = \sqrt{\frac{2(J(\boldsymbol{\theta}^0) - J(\boldsymbol{\theta}^*))}{L\sigma^2\bar{w}\bar{w}}} \quad (20)$$

in the above inequality, we get the desired result. \square

As seen in Theorem 1, our importance-weighted policy gradient estimator has a convergence rate of $O(1/\sqrt{K})$. Equivalently, the time complexity of our algorithm to obtain an ϵ -accurate solution is $O(1/\epsilon^2)$. Note that our choice of step size α requires knowing the total number of iterations K in advance. A more practical approach is to use a time-decayed step size of $\alpha_k \propto 1/\sqrt{k}$ or $\alpha_k \propto 1/k$.

B Case study

A demonstrative case study (Fig.9) comparing our proposed approach (LUFFY) against baseline methods (SFT and GPRO) in mathematical problem solving reveals distinct characteristics in reasoning patterns. SFT demonstrates redundant and circular reasoning with excessive repetition (over 8,129 tokens), while GPRO shows concise but unfounded deduction (1002 tokens), both leading to incorrect conclusions. In contrast, LUFFY presents a well-balanced approach (2623 tokens) that combines systematic decomposition with clear mathematical calculation. Through rigorous reasoning and proper verification steps, LUFFY successfully reaches the correct answer, demonstrating the effectiveness of our methodology in achieving both accuracy and efficiency.

Case Study

Prompt: <system prompt> This is the problem: A diagonal of a regular 2006-gon is called odd if its endpoints divide the boundary into two parts, each composed of an odd number of sides. Sides are also regarded as odd diagonals. Suppose the 2006-gon has been dissected into triangles by 2003 nonintersecting diagonals. Find the maximum possible number of isosceles triangles with two odd sides.

SFT

Response:

<think>

[Problem Restatement Not Analysis]

Okay, let's try to tackle this problem.

So, we have a regular 2006-gon, and we need to find the maximum number of isosceles triangles...

[Endless Definition Loop]

Let me break this down step by step. First, let's recall some basics about regular polygons and their diagonals...

Hmm, maybe they mean that sides are considered as diagonals here?...

Hmm, maybe they are using "diagonal" in a different sense here...

Alternatively, maybe the problem is using diagonal to mean any chord...

[Wait-Pattern Paralysis]

Wait, but in the problem statement...

Wait, but sides are chords...

Wait, but in the problem's definition, sides are odd diagonals...

[Small-Case Trap]

Let's take a regular polygon with a small even number of sides, say 4 (a square). A triangulation of a square has 2 triangles...

[Complete Breakdown]

but diagonal diagonal polygon. < odd diagonal in the odd of odd diagonal = think... wait diagonal diagonal diagonal diagonal...

Tokens Length: > 8192

Correctness: False

On-Policy RL

Response:

<think>

Formulaic Decomposition

To solve this problem, we need to understand ...

[Pseudo-Systematic Structure]

Let's step-by-step.

1. Understanding the ... if we label the vertices of the polygon as (V_1, ..., V_2006)...

[Circular Property Description]

Properties of Odd Diagonals... Properties of triangulation... Properties of the polygon...

[False Authority Appeal]

The number of isosceles triangles with two odd sides in a triangulation of a regular polygon is a well-known problem in combinatorial geometry...

[Unsubstantiated Conclusion]

For a 2006-gon, the maximum number of such triangles is 501.

</think>

[Final Answer]

The maximum possible number of isosceles triangles with two odd sides in a triangulation of a regular 2006-gon is

501.

Tokens Length: 1002

Correctness: False

LUFFY

Response:

<think>

Clear Problem Definition

To solve this problem, we need to understand ... and how it can be ... We are specifically interested in ...

[Systematic Decomposition and Analysis]

First, let's define some terms and properties...

Step 1: Understanding Odd Diagonals...

Step 2: Properties of Isosceles Triangles in a Regular Polygon...

Step 3: Counting Odd Diagonals...

Step 4: Dissecting the Polygon into Triangles...

[Precise Mathematical Calculation]

"For each vertex, there are 1002 vertices at an odd distance..."

[Logical Verification]

Given that each isosceles triangle with two odd sides uses one odd diagonal and one side...

[Consistency Check]

"Let's verify the logic:

1. Odd Diagonals ...

2. Dissection ..."

</think>

[Final Answer]

Thus, the maximum possible number of isosceles triangles with two odd sides is

1003.

Tokens Length: 2623

Correctness: True

Answer: "\$1003\$"

Figure 9: Comparison of three approaches(SFT, On-Policy RL, and LUFFY) for a geometric problem.