

# TOWARDS A UNIFIED VIEW OF LARGE LANGUAGE MODEL POST-TRAINING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Many approaches with seemingly disparate losses exist for post-training modern language models, such as Reinforcement Learning (RL) and Supervised Fine-Tuning (SFT). In this paper, we show that these approaches are not in contradiction, but are instances of a single optimization process. We derive the Unified Policy Gradient Estimator (UPGE), a framework with four interchangeable parts that unifies a wide spectrum of post-training approaches through their loss gradient form. We further present the calculations of these methods as the gradient of a common objective under different data distribution assumptions and various bias-variance tradeoffs. Motivated by our theoretical findings, we propose Hybrid Post-Training (HPT), an algorithm that dynamically selects different training signals. HPT is designed to yield both effective exploitation of demonstration and stable exploration without sacrificing learned reasoning patterns. We provide extensive experiments and ablation studies to verify effectiveness of HPT. Across six mathematical reasoning benchmarks and two out-of-distribution tasks, HPT consistently surpasses strong baselines across models of varying scales and families.

## 1 INTRODUCTION

Reinforcement Learning has played an integral role in enhancing the reasoning capabilities of large language models (LLMs) (Jaech et al., 2024; Team et al., 2025; Guo et al., 2025). However, applying RL directly to a base model (i.e., “Zero RL”) (Zeng et al., 2025b) presupposes a certain level of inherent capability. This method often falters when applied to weaker models or tasks of high complexity, as the exploration process may fail to discover meaningful reward signals. Conversely, the classical Supervised Fine-Tuning (SFT) (Wei et al., 2021) offers a direct and efficient method to distill knowledge from high-quality, human-annotated data, but often curtails the model’s exploratory capabilities. Consequently, a sequential “SFT-then-RL” pipeline (Yoshihara et al., 2025) has emerged as the standard, adopted by numerous state-of-the-art models. While effective, this multi-stage process is notoriously resource-intensive and usually requires careful tuning to ensure effectiveness.

To circumvent these challenges, recent works have focused on integrating SFT or SFT-style imitation learning losses directly with RL objectives (Yan et al., 2025; Fu et al., 2025; Zhang et al., 2025; Yu et al., 2025; Zeng et al., 2025a). In these approaches, the model is updated using a composite loss function. The balance between imitation and exploration components is governed by various strategies, including a fixed coefficient, a predefined schedule, a dynamic adjustment based on entropy, or a learnable parameter. Appendix A further discusses related work. These works predominantly treat SFT and RL losses as two distinct objectives. And a detailed analysis of *why these two learning signals can be effectively combined within a unified optimization process* remains largely unexplored.

Despite their distinct mathematical formulations, we find that the loss gradients of these approaches can be cast into a single form and jointly drive the optimization process under a unified framework. Inspired by Generalized Advantage Estimator (Schulman et al., 2015b), we introduce Unified Policy Gradient Estimator (UPGE), a framework that subsumes the gradients of various post-training objectives into one generalized expression with four interchangeable parts: *stabilization mask*, *reference policy*, *advantage estimate*, and *likelihood gradient*. We show that UPGE can be theoretically derived from a common objective, and the various forms of gradients are not conflicting but instead act as complementary learning signals that jointly guide optimization under this view. Further, through UPGE it becomes clear that these gradient estimators possess different characteristics, and there

Table 1: Theoretical unified view of various post-training algorithms of the large language model.

Algorithm	Reference Policy	Advantage Estimate	Unified Policy Gradient Estimator
SFT	$\pi_{ref} = \pi_\theta$	$\hat{A}_{SFT} \equiv 1$	$\nabla \mathcal{J}_{SFT}(\theta) = \nabla \pi_\theta(\tau) \frac{\hat{A}_{SFT} \equiv 1}{\pi_\theta(\tau)}$
<b>Online Reinforcement Learning Methods</b>			
PPO (Schulman et al., 2017)	$\pi_{ref} = \pi_{\theta_{old}}$	$\hat{A}_{PPO} = \text{GAE}$ (Schulman et al., 2015b)	$\nabla \mathcal{J}_{PPO} = \nabla \pi_\theta(\tau) \frac{\hat{A}_{PPO} \hat{A}_{Clip}}{\pi_{ref}(\tau)}$
GRPO (Shao et al., 2024)	$\pi_{ref} = \pi_{\theta_{old}}$	$\hat{A}_{GRPO} = \frac{R(\tau_j) - \text{mean}(\{R(\tau_j)\}_{G_{on}})}{\text{std}(\{R(\tau_j)\}_{G_{on}})}$	$\nabla \mathcal{J}_{GRPO} = \nabla \pi_\theta(\tau) \frac{\hat{A}_{GRPO} \hat{A}_{Clip}}{\pi_{ref}(\tau)}$
REINFORCE (Ahmadian et al., 2024)	$\pi_{ref} = \pi_\theta$	$\hat{A}_{REINFORCE} = \pm 1$	$\nabla \mathcal{J}_{REF}(\theta) = \nabla \pi_\theta(\tau) \frac{\hat{A}_{REF}}{\pi_\theta(\tau)}$
CISPO (Chen et al., 2025)	$\pi_{ref} = \pi_{\theta_{old}}$	$\hat{A}_{CISPO} = \hat{A}_{GRPO}$	$\nabla \mathcal{J}_{CISPO} = \nabla \pi_\theta(\tau) \frac{\hat{A}_{CISPO} \hat{A}_{Clip} \text{Mask}}{\pi_{ref}(\tau)}$
GSPO (Zheng et al., 2025)	$\pi_{ref} = \pi_\theta \left( \frac{\pi_{\theta_{old}}(\tau_{i,j} q_i)}{\pi_\theta(\tau_{i,j} q_i)} \right)^{1/ \tau_{i,j} }$	$\hat{A}_{GSPO} = \hat{A}_{GRPO}$	$\nabla \mathcal{J}_{GSPO} = \nabla \pi_\theta(\tau) \frac{\hat{A}_{GSPO} \hat{A}_{Clip}}{\pi_{ref}(\tau)}$
<b>Offline/Online Reinforcement Learning Methods</b>			
SRFT (Offline) (Fu et al., 2025)	$\pi_{ref} \equiv 1$	$\hat{A}_{SRFT} = \frac{R(\tau_j) - \text{mean}(\{R(\tau_j)\}_{G_{on} \cup G_{off}})}{\text{std}(\{R(\tau_j)\}_{G_{on} \cup G_{off}})}$	$\nabla \mathcal{J}_{SRFT} = \nabla \pi_\theta(\tau) \frac{\hat{A}_{SRFT}}{\pi_{ref}(\tau) \equiv 1}$
LUFFY (Offline) (Yan et al., 2025)	$\pi_{ref} \equiv 1$	$\hat{A}_{LUFFY} = \hat{A}_{SRFT}$	$\nabla \mathcal{J}_{LUFFY} = \nabla \pi_\theta(\tau) \frac{\hat{A}_{LUFFY} f'_{shape}}{\pi_{ref}(\tau) \equiv 1}$

exists a bias–variance tradeoff in their respective gradient components, which suggests that an ideal post-training algorithm should dynamically select among them to suit different conditions. Building upon this unified perspective, we propose Hybrid Post-Training (HPT), a hybrid algorithm to dynamically choose more desirable training signals by switching between SFT and RL based on rollout accuracy feedback. This mechanism allows HPT to be intrinsically adaptive to models of varying capabilities and data of differing complexities.

Empirical evaluations demonstrate that HPT surpasses baselines such as SFT→GRPO and LUFFY on Qwen2.5-Math-7B, achieving a 7-point gain over our strongest baseline on AIME 24. Moreover, HPT yields substantial improvements even on smaller and weaker models, including Llama3.1-8B and Qwen2.5-Math-1.5B. Through detailed training dynamics and illustrative visualizations, we clearly reveal the features and underlying mechanisms of HPT. The following are several key takeaways:

**Takeaways**

1. UPGE provides a theoretical unification of a wide spectrum of post-training algorithms, covering both SFT and RL losses within a single formulation (§ 2).
2. HPT is capable of outperforming previous post-training and mixed-policy algorithms across a diverse range of models and benchmarks (§ 3).
3. Dynamic integration of SFT and RL in HPT achieves the highest *Pass@1024*, facilitating enhanced exploration and generalization of the model (§ 4.1).

## 2 A UNIFIED VIEW ON POST-TRAINING ALGORITHMS

Despite the substantial differences in loss formulations across various LLM post-training algorithms, we find that their loss gradients can be expressed in a single common form. We propose a unified framework for the gradient calculation, named **Unified Policy Gradient Estimator (UPGE)**:

$$\text{grad}_{Uni} = \mathbb{1}_{stable} \frac{1}{\pi_{ref}} \hat{A} \nabla \pi_\theta.$$

In the following sections, we first provide a detailed analysis of the components of UPGE (§ 2.1). Then we theoretically derive the UPGE from a common objective (§ 2.2). Based on the unified perspective, we finally propose the Hybrid Post-Training (HPT) algorithm (§ 2.3).

### 2.1 COMPONENTS OF THE UNIFIED POLICY GRADIENT ESTIMATOR

We present the Unified Policy Gradient Estimator, our unified framework for gradient calculations. In Table 1, we list the UPGE instantiations for some representative post-training methods, with detailed

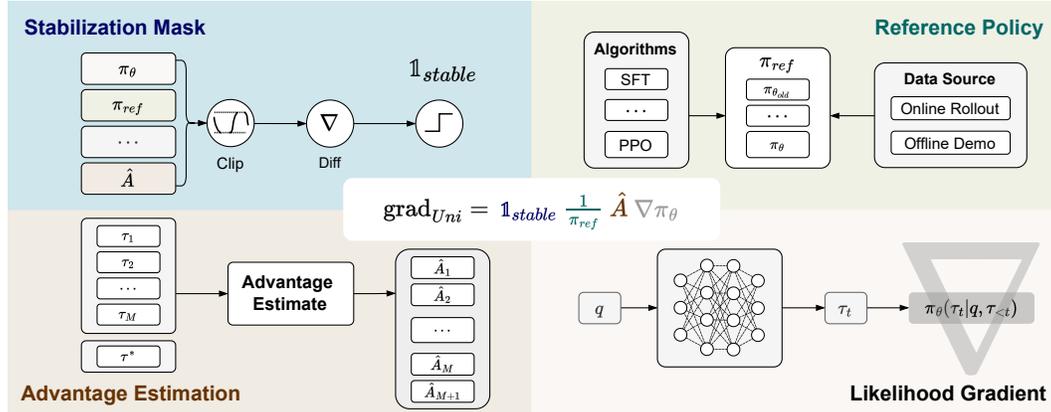


Figure 1: Illustration of the Unified Policy Gradient Estimator. The “ $\nabla$ ” in the background of the Likelihood Gradient part refers to the calculation of the gradient with respect to the  $\pi_\theta$ .

derivations provided in Appendix B. UPGE comprises four components: *stabilization mask*, *reference policy*, *advantage estimate*, and *likelihood gradient*. We address each of the terms below.

**Stabilization Mask**  $\mathbb{1}_{stable}$  Starting from PPO (Schulman et al., 2017), the stabilization mask was first derived as an approximation of the TRPO Algorithm (Schulman et al., 2015a). In practice, the PPO clipping addresses the instability issue during RL training by turning off the current gradient when the current iterate is considered unsafe. In consequent works in Table 1, many have provided their modifications on the stability mask, usually motivated by empirical evaluations.

**Reference Policy Denominator**  $\pi_{ref}$  We note that our notion of reference policy on the denominator differs from the commonly used rollout policy  $\pi_{old}$ . This denominator denotes a token-level reweight coefficient, usually in the form of an inverse probability. There are multiple choices for this coefficient. For the case of SFT, the policy denominator uses the current policy  $\pi_\theta(\tau)$ , which is a result of  $\mathcal{L} = -\log(\pi_\theta(\tau))$  as the objective function, where  $\tau$  denotes a trajectory. For the case of PPO-style online RL algorithms, generally, the policy denominator uses the rollout policy  $\pi_{old}(\tau)$ .

**Advantage Estimate**  $\hat{A}$  In traditional RL, the advantage evaluates the additional benefit of taking the current action given the state. Similarly, the post-training process seeks to maximize the likelihood of generating positive sequences with high advantage and minimize negative sequences.

**Likelihood Gradient**  $\nabla\pi_\theta(\tau)$  Likelihood gradient is a general term which maps gradient information from the actions to the model parameters  $\theta$ . It is crucial for back-propagating the objective signals to the network weights, and is kept the same formal structure across all gradient calculations.

Within UPGE, existing post-training algorithms can be analysed under these four components respectively, clarifying their distinctive properties, which is presented in detail in Appendix C.

## 2.2 DERIVATION OF THE UNIFIED POLICY GRADIENT ESTIMATOR

In this section, we theoretically derive the UPGE from a common objective shared by all post-training algorithms: improve the likelihood of positive trajectories and decrease that of negative trajectories such that the total reward in expectation  $\max_\theta \mathcal{J}(\theta) := \mathbb{E}[r(\tau|q)]$  is maximized. We then show that SFT and RL objectives are not in conflict, and they can be optimized jointly within a single loss.

**Common Objective.** We model the post-training as a process to maximize the expected success rate while keeping the model policy closely adhering to a demonstration dataset (behavior policy)  $\pi_\beta$ :

$$\mathcal{J}_\mu(\theta) = \mathbb{E}_{\tau \sim \pi_\theta(\cdot|q)}[r(\tau|q)] - \mu \text{KL}(\pi_\beta(\cdot|q) \parallel \pi_\theta(\cdot|q)), \quad \mu \geq 0, \quad (1)$$

where  $q \sim \mathcal{D}$  denotes the question from a given distribution,  $\tau$  denotes a trajectory,  $r$  denotes the (binary/real) score, and  $\pi_\beta$  denotes behavior policy from demonstration.

**Gradient of Common Objective.** Differentiating Eq.1 (full derivation in Appendix D.1), we obtain

$$\nabla_{\theta} \mathcal{J}_{\mu}(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [r(\tau | q) \nabla_{\theta} \log \pi_{\theta}(\tau | q)] + \mu \mathbb{E}_{\tau \sim \pi_{\beta}} [\nabla_{\theta} \log \pi_{\theta}(\tau | q)]. \quad (2)$$

**From gradient to UPGE.** Applying measure-change identity (detailed in Appendix D.1) with reference policy  $\pi_{ref}$  mentioned in Section 2.1 and using  $\nabla \log \pi_{\theta} = (1/\pi_{\theta}) \nabla \pi_{\theta}$  yields the gradient:

$$\nabla_{\theta} \mathcal{J}_{\mu}(\theta) = \mathbb{E}_{\tau \sim \pi_{ref}(\cdot | q)} \left[ \frac{1}{\pi_{ref}(\tau | q)} \widehat{A}_{uni}(\tau, q) \nabla_{\theta} \pi_{\theta}(\tau | q) \right], \quad (3)$$

$$\widehat{A}_{uni}(\tau, q) = \underbrace{r(\tau | q)}_{\widehat{A}_{RL}(\tau, q)} + \underbrace{\mu \mathbb{1}\{\pi_{ref} = \pi_{\beta}\} \frac{\pi_{\beta}(\tau | q)}{\pi_{\theta}(\tau | q)}}_{\widehat{A}_{SFT}(\tau, q)}. \quad (4)$$

In many RL works, the raw score  $r(\tau | q)$  is replaced by a more structured advantage to reduce variance, provide relative credit assignment within a rollout group, and stabilize step sizes. For example, GRPO uses group-wise normalization:  $\widehat{A}_{GRPO}(\tau_j, q) = \frac{R(\tau_j) - \text{mean}(\{R(\tau)\}_{G_{on}})}{\text{std}(\{R(\tau)\}_{G_{on}})}$ .

When trust-region stabilization masks, as induced by PPO clipping (detailed in Appendix D.3), are inserted multiplicatively without altering the target objective, we obtain our UPGE:

$$\text{grad}_{uni} = \mathbb{E}_{\tau \sim \pi_{ref}(\cdot | q)} \left[ \mathbb{1}_{stable}(\tau, q) \frac{1}{\pi_{ref}(\tau | q)} \widehat{A}_{uni}(\tau, q) \nabla_{\theta} \pi_{\theta}(\tau | q) \right] \quad (5)$$

Details of the derivation above are provided in Appendix D. The gradient in Eq. 2 is the sum of a reward term sampled from  $\pi_{\theta}$  and a data-adherence (SFT) term sampled from  $\pi_{\beta}$ . Both terms map to the same estimator via Eq. 3–5 by choosing  $\pi_{ref}$  accordingly. (e.g.,  $\pi_{\theta_{old}}$  for on-policy trust-region updates and  $\pi_{\beta}$  for SFT/offline updates). Therefore, SFT and RL optimize a single Common Objective (Eq. 1) and can be trained jointly within one loss without intrinsic conflict. While these algorithms share the Common Objective, they retain distinct characteristics, and bias-variance trade-offs still exist for different components of the unified gradient estimator, as shown in Appendix C. This observation suggests that an ideal post-training method should adaptively select among methods under varying conditions.

### 2.3 HYBRID POST-TRAINING WITH PERFORMANCE FEEDBACK

Our unified perspective above shows that different post-training losses have the same optimization objective with different characteristics. Inspired by this view, we propose the Hybrid Post-Training (HPT) algorithm. We use a mixed loss  $\mathcal{L} = \alpha \mathcal{L}_{RL} + \beta \mathcal{L}_{SFT}$ , which contains the weighted on-policy RL loss  $\mathcal{L}_{RL}$  and SFT loss  $\mathcal{L}_{SFT}$ , to optimize the target LLM  $\pi_{\theta}$ . The weights of the two losses ( $\alpha$  and  $\beta$ ) are determined by the real-time sampling performance of the model.

**Performance on Single Question.** For each question  $q$  provided to LLM, we first obtain supervising trajectory  $\tau^*$  and model’s performance  $P$  on  $q$ . Specifically, we draw  $n$  on-policy trajectories  $\{\tau_i\}_{i=1}^n \sim \pi_{\theta}(\cdot | q)$  and evaluate them with a verifier  $v : \tau_i \rightarrow \{0, 1\}$ , which is the same as the rule-based reward function. Performance  $P$  is defined as the mean of these  $n$  verification scores:

$$P = \frac{1}{n} \sum_{i=1}^n v(\tau_i), \quad v(\tau_i) = R(\tau_i) = \begin{cases} 1 & \text{if } \tau_i \text{ contains the correct answer of } q \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

**Feedback Coefficients.** Then, we obtain the coefficients based on the performance feedback:

$$\alpha = f(P), \quad \beta = g(P), \quad (7)$$

where the  $f$  and  $g$  are the specific feedback functions. Experientially, when the model demonstrates strong capability, it is advantageous to emphasize on-policy RL to foster exploration; conversely, when the model’s competence is limited, SFT should take precedence to ensure correct guidance. Consequently,  $f$  ought to increase with  $P$ , while  $g$  should decrease. In this paper, we employ a pair of simple yet empirically effective switch functions  $f$  and  $g$ :

$$\alpha = f(P) = \begin{cases} 1 & \text{if } P > \gamma \\ 0 & \text{if } P \leq \gamma \end{cases}, \quad \beta = g(P) = \begin{cases} 1 & \text{if } P \leq \gamma \\ 0 & \text{if } P > \gamma \end{cases} \quad (8)$$

The switch gate  $\gamma$  enables the model to perform SFT when its performance falls below a predefined threshold, and RL otherwise.

**Mixed Loss.** Finally, we calculate the RL loss  $\mathcal{L}_{\text{RL}}$  using Dr. GRPO with the already generated  $n$  on-policy trajectories  $\tau_i$  and SFT loss  $\mathcal{L}_{\text{SFT}}$  with the supervising trajectory  $\tau^*$ :

$$\mathcal{L}_{\text{RL}} = -\frac{1}{n} \sum_{i=1}^n \sum_{t=1}^{|\tau_i|} \min\left(r_{i,t} A_{i,t}, \text{clip}(r_{i,t}, 1 - \epsilon, 1 + \epsilon) A_{i,t}\right) \quad (9)$$

$$\mathcal{L}_{\text{SFT}} = -\frac{1}{|\tau^*|} \sum_{t=1}^{|\tau^*|} \log \pi_{\theta}(\tau_t^* | q, \tau_{<t}^*) \quad (10)$$

where  $r_{i,t} = \frac{\pi_{\theta}(\tau_{i,t} | q, \tau_{i,<t})}{\pi_{\theta_{\text{old}}}(\tau_{i,t} | q, \tau_{i,<t})}$  is the per-token importance sampling ratio,  $A_{i,t} \equiv A_i = \frac{R(\tau_i) - \text{mean}(\{R(\tau_i) \mid i=1,2,\dots,n\})}{\text{std}(\{R(\tau_i) \mid i=1,2,\dots,n\})}$  is advantage and  $\epsilon$  is the clip gate. The mixed loss is then obtained by taking a weighted average of these two losses using performance feedback coefficients  $\alpha$  and  $\beta$ :

$$\mathcal{L} = \alpha \mathcal{L}_{\text{RL}} + \beta \mathcal{L}_{\text{SFT}} \quad (11)$$

The algorithm procedure and detailed analysis of the effectiveness and advantages of our HPT are provided in Appendix E.

### 3 MAIN EXPERIMENTS

#### 3.1 EXPERIMENTAL SETUP

**Models** To evaluate the generalizability of HPT across different models, we conduct experiments using Qwen and LLaMA models of various scales. The models we experiment with are Qwen2.5-Math-1.5B, Qwen2.5-Math-7B (Yang et al., 2024), and LLaMA-3.1-8B (Grattafiori et al., 2024).

**Evaluation Setup** We evaluate HPT on 6 mathematical reasoning benchmarks: AIME 2024 (Li et al., 2024), AIME 2025 (Li et al., 2024), AMC (Li et al., 2024), MATH-500 (Hendrycks et al., 2021a), Minerva (Lewkowycz et al., 2022), and OlympiadBench (He et al., 2024). Moreover, when employing Qwen2.5-Math-7B as the backbone, we further conduct evaluations on ARC-c (Clark et al., 2018) and GPQA-Diamond (Rein et al., 2024). For the main experiments, following DeepSeek-R1 (Guo et al., 2025), we adopt the *Pass@k* evaluation protocol (Chen et al., 2021) and report *Pass@1* using non-zero temperature sampling. To ensure a fair comparison with previous works (Yan et al., 2025; Fu et al., 2025), we compute *avg@32* for AIME 24, AIME 25, and AMC (*avg@1* for others) using a temperature of 0.6 and a top- $p$  value of 0.95 for accuracy calculation.

**Baselines** Since HPT dynamically integrates GRPO (Shao et al., 2024) and SFT, the most natural baselines are SFT and GRPO individually. Furthermore, we compare HPT against the mix-policy approach LUFFY (Yan et al., 2025). For experiments using Qwen2.5-Math-7B as the backbone, we additionally include SFT→GRPO and SRFT<sup>1</sup> (Fu et al., 2025) as a baseline, as well as models trained with the Zero-RL procedure on the same backbone for a more comprehensive comparison. We also use PRIME-Zero (Cui et al., 2025a), SimpleRL-Zero (Zeng et al., 2025c), OpenReasoner-Zero (Hu et al., 2025b) and Oat-Zero (Liu et al., 2025b) as baselines.

**Implementation Details** We fix the switch gate  $\gamma$  at 0 throughout all experiments on the Qwen Family models and 2 for LLaMA, and provide relative ablation studies in Section G.2. For hyperparameters, we use a constant learning rate of  $5 \times 10^{-6}$  and adopt the AdamW optimizer for the policy model. For rollout, we sample 8 responses using a temperature of 1.0. The maximum generation length is set to 8, 192 tokens for all models. Other details are reported in Appendix F.

#### 3.2 MAIN RESULTS

Table 2 presents the overall performance of HPT on Qwen2.5-Math-7B. HPT not only significantly outperforms both SFT-only and GRPO-only baselines, but also surpasses SFT→GRPO, which requires substantially higher computational cost. This suggests that simply concatenating the two

<sup>1</sup>The results of SRFT are based on our own implementation, as the official code is not public.

Table 2: In-distribution and out-of-distribution performance of our HPT and baselines on the Qwen2.5-Math-7B. \* means the results are taken from the corresponding paper.

Model	In-Distribution							Out-of-Distribution		
	AIME 24	AIME 25	AMC	MATH-500	Minerva	Olympiad	Avg	ARC-c	GPQA	Avg
Qwen2.5-Math-7B	12.3	4.7	33.0	43.6	8.8	13.6	19.3	30.9	28.3	29.6
SFT	25.1	<b>22.8</b>	56.1	84.2	33.8	44.7	44.5	67.4	25.3	46.4
GRPO	19.4	13.8	59.1	81.8	38.2	46.2	43.1	81.2	36.4	58.8
SFT → GRPO	25.7	21.6	62.2	84.6	38.2	46.8	46.5	67.7	30.8	49.3
LUFFY	26.1	21.8	66.2	88.4	41.9	54.1	49.8	80.8	39.4	60.1
SRFT	18.4	15.5	55.9	83.8	42.6	48.9	44.2	80.5	36.8	58.7
HPT	<b>33.0</b>	21.9	<b>69.4</b>	<b>89.2</b>	<b>46.0</b>	<b>56.9</b>	<b>52.7</b>	<b>81.6</b>	<b>42.9</b>	<b>62.3</b>
Qwen2.5-Math-7B-Ins.	11.8	9.8	48.3	83.2	34.2	39.3	37.8	72.7	29.3	51.0
PRIME-Zero*	17.0	12.8	54.0	81.4	39.0	40.3	40.8	73.3	18.2	45.8
SimpleRL-Zero*	27.0	6.8	54.9	76.0	25.0	34.7	37.4	30.2	23.2	26.7
OpenReasoner-Zero*	16.5	15.0	52.1	82.4	33.1	47.1	41.0	66.2	29.8	48.0
Oat-Zero*	33.4	11.9	61.2	78.0	34.6	43.4	43.8	70.1	23.7	46.9

training stages is not the most effective strategy. Moreover, HPT achieves marked improvements over existing mixed-policy approaches such as LUFFY and SRFT, with particularly notable gains of 6.9 and 14.6 points on AIME 2024, respectively. Furthermore, we conduct experiments on models of different scales and families to evaluate the effectiveness of HPT, including LLaMA3.1-8B and Qwen2.5-Math-1.5B, as shown in Table 3. Compared with SFT, GRPO, and LUFFY, HPT achieves substantial performance gains, demonstrating robustness across models with varying capability levels.

Table 3: Performance of HPT and baselines on LaMA3.1-8B and Qwen2.5-Math-1.5B. \* means the results are taken from the LUFFY paper (Yan et al., 2025).

Model	AIME 24	AIME 25	AMC	MATH-500	Minerva	Olympiad	Avg
LLaMA3.1-8B	0.4	0.1	4.7	13.8	4.8	3.9	4.6
SFT*	0.5	0.1	5.4	20.2	4.0	5.3	5.9
GRPO*	0.3	0.5	9.4	23.4	17.6	6.1	9.6
LUFFY*	1.9	0.1	13.5	39.0	15.1	9.6	13.2
HPT	<b>2.1</b>	<b>1.2</b>	<b>18.6</b>	<b>47.8</b>	<b>18.8</b>	<b>20.4</b>	<b>18.2</b>
Qwen2.5-Math-1.5B	2.8	6.1	24.5	32.8	11.0	16.4	15.6
SFT	14.7	17.6	45.4	78.4	29.4	35.7	36.9
GRPO	12.2	8.5	43.8	71.0	33.1	35.3	34.0
LUFFY	14.1	9.4	43.5	75.2	26.1	39.7	34.7
HPT	<b>16.6</b>	<b>17.8</b>	<b>51.0</b>	<b>81.0</b>	<b>37.5</b>	<b>47.3</b>	<b>41.9</b>

## 4 EMPIRICAL ANALYSIS

Our empirical analysis progressively reveals how HPT reconciles exploration and exploitation, stabilizes training, and ultimately enhances the reasoning ability. We begin in § 4.1 with an examination of *exploration* and *exploitation*. In § 4.2, we provide a training visualization, contrasting HPT with the conventional SFT→GRPO. Finally, § 4.3 investigates fine-grained training metrics of HPT.

### 4.1 EXPLORATION AND EXPLOITATION

HPT inherently achieves an adaptive switching between RL and SFT. These two paradigms naturally correspond to the learning modes of *exploration* and *exploitation*. Accordingly, we can examine whether HPT addresses the initial challenges from both perspectives.

**Exploration** From the exploration perspective, we want to analyze the model’s  $Pass@k$  performance after training with HPT. Recently, Limit-of-RLVR (Yue et al., 2025) demonstrated that while RLVR yields an improvement in  $Pass@1$ , it does not lead to gains in large- $k$   $Pass@k$ . In other words, RLVR does not expand the exploratory capability boundary of the base model. We follow Yue et al. (2025) to evaluate  $Pass@k$  up to 1024 for each question of AIME25, AIME24, and AMC for  $Pass@k$  evaluation. Based on these sets of generated solutions, we apply bootstrap sampling to obtain accurate estimates of  $Pass@k$  scores for various values of  $k$ . Figure 2 illustrates the resulting  $Pass@k$  curves.

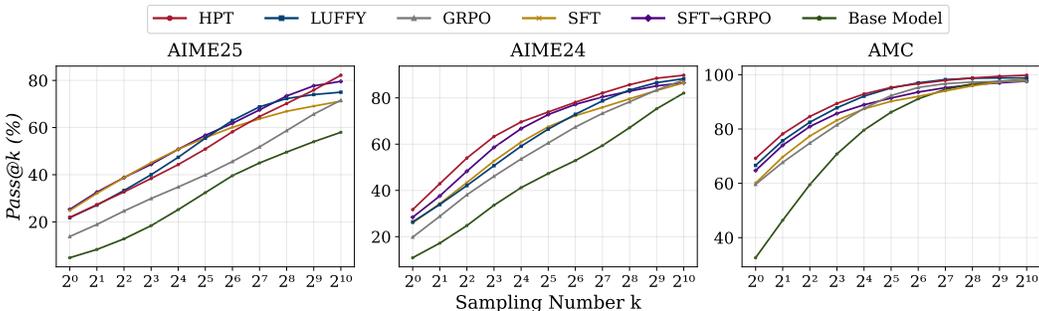


Figure 2:  $Pass@k$  performance of HPT against baselines on Qwen2.5-Math-7B. Evaluation spans 3 benchmarks, with  $Pass@k$  values estimated via bootstrap sampling from a set of 2048 generated solutions per question.

- First, we can observe that methods incorporating SFT achieve higher large- $k$   $Pass@k$  compared to the GRPO (purely RL). This may be attributed to the introduction of data outside the model’s own distribution during SFT, which increases output uncertainty while also providing new knowledge from offline data, thereby enhancing the model’s exploratory capacity.
- Furthermore, we identify an interesting phenomenon: since HPT dynamically integrates RL (GRPO) with SFT, we might intuitively expect its large- $k$   $Pass@k$  performance to fall between that of the two individual methods. However, HPT achieves the highest large- $k$   $Pass@k$  performance overall. This indicates that *Hybrid Post-Training not only delivers substantial improvements in  $Pass@1$ , but also maximally preserves and enhances the model’s exploratory ability.*

Table 4: Bidirectional analysis of exclusive solves on MATH-500 with Qwen2.5-Math-7B as the backbone model, comparing the HPT against baselines (GRPO and LUFFY). The notation  $+X/-Y$  in each cell indicates the performance trade-off:  $+X$  represents the number of problems solved by the HPT but not the baseline, while  $-Y$  represents the number solved by the baseline but not by the HPT.

Methods	Level 1 (N=43)	Level 2 (N=90)	Level 3 (N=105)	Level 4 (N=128)	Level 5 (N=134)	Overall (N=500)
<b>GRPO</b>						
Absolute	+0/-0	+5/-1	+9/-2	+17/-4	+27/-8	+58/-15
Percentage	+0.0%/-0.0%	+5.6%/-1.1%	+8.6%/-1.9%	+13.3%/-3.1%	+20.1%/-6.0%	+11.6%/-3.0%
<b>LUFFY</b>						
Absolute	+1/-0	+5/-1	+5/-3	+10/-5	+22/-7	+43/-16
Percentage	+2.3%/-0.0%	+5.6%/-1.1%	+4.8%/-2.9%	+7.8%/-3.9%	+16.4%/-5.2%	+8.6%/-3.2%

**Exploitation** From the exploitation perspective, the key question is whether HPT, by leveraging SFT, enhances the model’s initial competence and facilitates subsequent RL training. To investigate this, we analyze its exclusive solves against the GRPO and LUFFY, building upon the results from the evaluation on MATH-500, as shown in Table 4. The red numbers denote problems that are solved by our HPT but not by GRPO or LUFFY, i.e., problems newly acquired through our training procedure.

- First, the red counts consistently increase with problem difficulty, suggesting that HPT improves the model’s ability to tackle more challenging problems. And the red counts are consistently large relative to both baselines, demonstrating that our method enables the model to acquire a substantial number of problems that prior approaches struggled to solve.
- Furthermore, the green counts remain essentially unchanged across settings, which indicates that, compared with existing methods, HPT preserves performance on problems that the model could already solve, thereby mitigating the risk of catastrophic forgetting.

#### 4.2 TRAINING VISUALIZATION

To facilitate a fine-grained examination of the training process and obtain deeper insights into how HPT works, we conduct a visualization analysis comparing the SFT→GRPO with HPT. We sample 255 questions from the MATH dataset (Hendrycks et al., 2021b) for training, with 85 each from Levels 3, 4, and 5. For SFT→GRPO, we perform 50 epochs of GRPO on a Qwen2.5-Math-1.5B

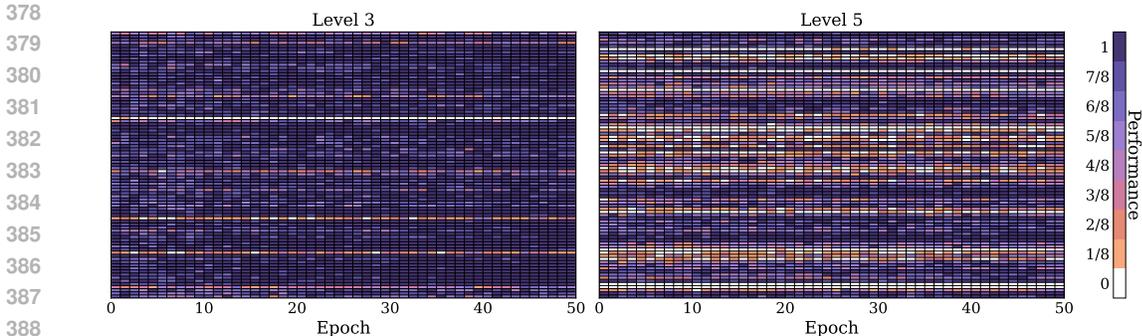


Figure 3: GRPO training dynamics of SFT→GRPO on Qwen2.5-Math-1.5B across 50 epochs. Each line represents a question and we visualize the model’s per-question rollout accuracy throughout the training process.

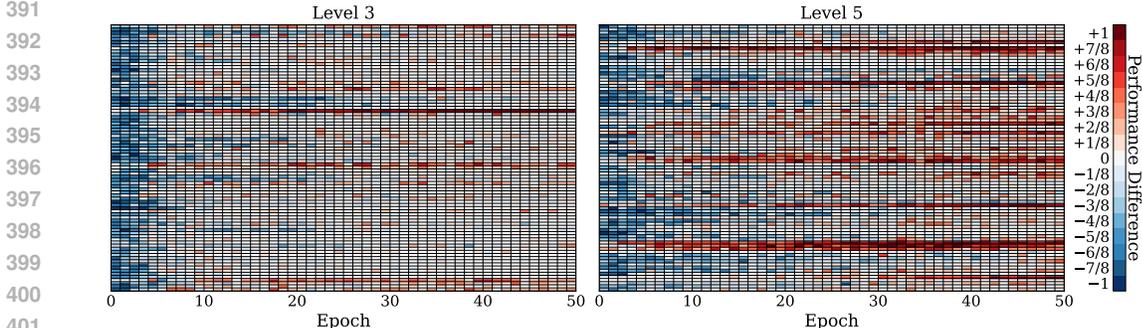


Figure 4: Performance difference (HPT v.s. SFT→GRPO) on Qwen2.5-Math-1.5B across 50 epochs. Each line represents a question. A diverging color scale indicates the advantage: red for HPT, blue for SFT→GRPO.

model fine-tuned with SFT, tracking rollout accuracy of each question across training, as shown in Figure 3. To highlight difficulty effects, we focus on Levels 3 and 5 as representative cases. Notably, GRPO frequently produces continuous white lines, reflecting widespread rollout errors across outputs. This illustrates a core limitation of RL methods: they struggle to learn effectively when frequent rollout errors occur.

In parallel, we use HPT to train Qwen2.5-Math-1.5B from scratch. To enable a more intuitive comparison, we calculate the accuracy difference at matched questions and epochs in the evaluation grid between two methods: red indicates HPT is better, blue the opposite. Figure 4 presents the results of the difference plots. Notably, SFT→GRPO, which involves a preceding SFT phase, requires greater computational resources than HPT. This additional computation leads to an initial dominance of the blue regions, considering the SFT stage in SFT→GRPO has already incorporated substantial prior knowledge. However, in the later training stage, HPT still surpasses and ultimately reveals the dominance of the red regions, indicating its superiority. This advantage becomes even more pronounced in the Level 5 subplot, suggesting that HPT provides particular benefits for learning on more challenging problems, which may be due to its appropriate use of the demonstration data.

### 4.3 TRAINING DYNAMICS

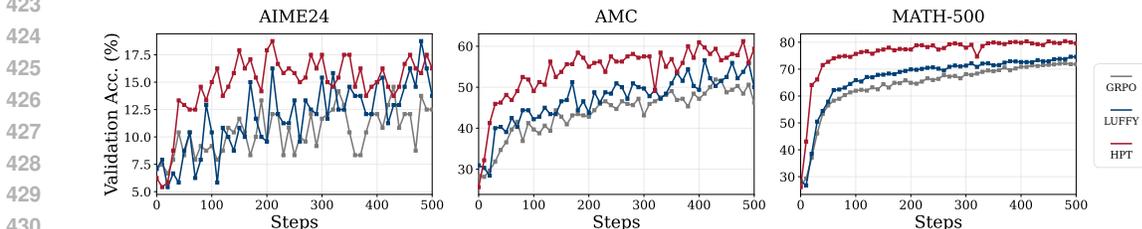


Figure 5: Validation performance comparisons on Qwen2.5-Math-1.5B across AIME24, AMC and MATH-500.

**Validation Performance.** We track the validation performance on Qwen2.5-Math-1.5B as shown in Figure 5, where HPT consistently outperforms the baselines and delivers stable improvements.

**Offline Data Ratio.** As shown in Figure 6, we quantify the offline data ratio, which is defined as the proportion of offline samples relative to the total number of training samples, at each training step. As expected, when the model has not yet acquired competence on the target tasks, the early phase is characterized by a large proportion of SFT-driven updates. As training progresses and the model’s on-policy reward increases, the mixture gradually shifts: the contribution of RL grows while that of SFT diminishes, eventually stabilizing at a small but non-zero level. This trend is observed for both Qwen2.5-Math-7B and Qwen2.5-Math-1.5B. The weaker 1.5B model remains in the SFT-dominated regime for a longer period before transitioning, whereas the stronger 7B model shifts earlier. These results align with our analysis of the design of HPT, where the mixing ratio is automatically adjusted based on performance rather than fixed in advance like LUFFY.

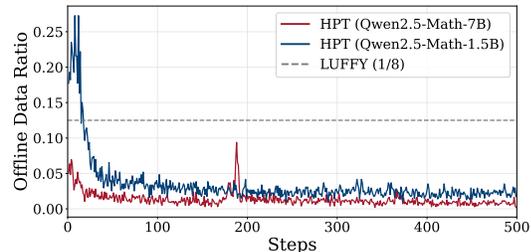


Figure 6: Offline data ratio dynamics during training, which is calculated as the proportion of offline training samples relative to the total training data at each step.

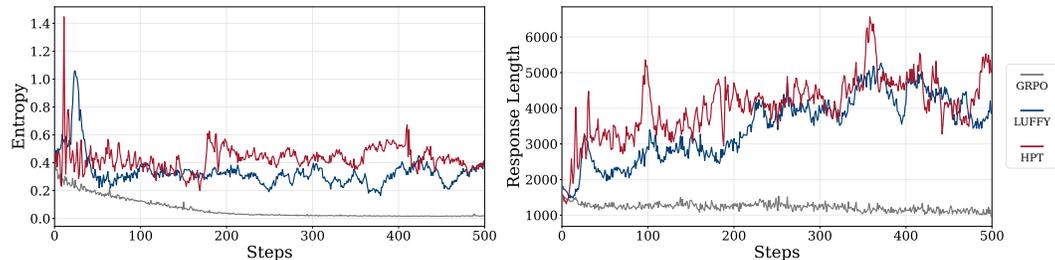


Figure 7: Comparisons of training dynamics: (left) The entropy measures the diversity of model outputs, indicating exploration behavior. (right) The response length tracks the average length of generated responses.

**Entropy and exploration.** Figure 7 (left) tracks token-level entropy over 500 steps. HPT maintains higher entropy than GRPO throughout the training phases. This is expected as the offline SFT trajectories are derived from the external demonstration distribution, which consequently increases the diversity in the model’s outputs.

**Response length and acquired reasoning patterns.** Figure 7 (right) reports the average response length. Our SFT trajectories have a length of up to 8k tokens. Under HPT, the model’s response length increases quickly during the early steps but does not jump to the 8k ceiling. More importantly, after the method shifts toward RL and the SFT proportion plateaus at a low level, the response length does not regress. This persistence suggests that the model has internalized long-form reasoning routines from the offline data rather than merely echoing teacher outputs. In other words, the learned reasoning pattern becomes part of the policy, and RL fine-tuning refines it instead of erasing it.

## 5 CONCLUSION

In this paper, we introduce the Unified Policy Gradient Estimator to unify the post-training of LLMs. We demonstrate that SFT and RL optimize a common objective, with their gradients representing different bias-variance tradeoffs. Motivated by this unified perspective, we propose Hybrid Post-Training (HPT), an algorithm that dynamically adapts between SFT for exploitation and RL for exploration based on real-time performance feedback. Extensive empirical validation shows that HPT consistently outperforms strong baselines across various models and benchmarks. Our work contributes both a unifying theoretical perspective on post-training and a practical algorithm that effectively balances exploitation and exploration to enhance model capabilities.

486 ETHICS STATEMENT  
487

488 This paper introduces UPGE to unify the post-training of LLMs and the mix-policy algorithm HPT.  
489 We study post-training for LLM reasoning on public math/logic benchmarks only, and no human  
490 subjects or personal data are used. We think our work doesn't introduce any ethical concerns.

491 REPRODUCIBILITY STATEMENT  
492

493 We provide the implementation code and scripts of our HPT algorithm and main baselines in the  
494 Supplementary Material for reproduction. And we also specify the experimental configurations and  
495 hyperparameters in detail in Section 3 and Appendix F.

496 REFERENCES  
497

- 498 Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin,  
499 Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning  
500 from human feedback in llms. *arXiv preprint arXiv:2402.14740*, 2024.
- 501 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain,  
502 Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with  
503 reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- 504 Aili Chen, Aonian Li, Bangwei Gong, Binyang Jiang, Bo Fei, Bo Yang, Boji Shan, Changqing Yu,  
505 Chao Wang, Cheng Zhu, et al. Minimax-ml: Scaling test-time compute efficiently with lightning  
506 attention. *arXiv preprint arXiv:2506.13585*, 2025.
- 507 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared  
508 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large  
509 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 510 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
511 reinforcement learning from human preferences. In *Advances in neural information processing  
512 systems*, volume 30, 2017.
- 513 Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V  
514 Le, Sergey Levine, and Yi Ma. Sft memorizes, rl generalizes: A comparative study of foundation  
515 model post-training. *arXiv preprint arXiv:2501.17161*, 2025.
- 516 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li,  
517 Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language  
518 models. *arXiv preprint arXiv:2210.11416*, 2022.
- 519 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and  
520 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.  
521 *arXiv preprint arXiv:1803.05457*, 2018.
- 522 Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang, Wendi Li, Bingxiang He, Yuchen Fan, Tianyu  
523 Yu, Qixin Xu, Weize Chen, Jiarui Yuan, Huayu Chen, Kaiyan Zhang, Xingtai Lv, Shuo Wang,  
524 Yuan Yao, Xu Han, Hao Peng, Yu Cheng, Zhiyuan Liu, Maosong Sun, Bowen Zhou, and Ning  
525 Ding. Process reinforcement through implicit rewards, 2025a.
- 526 Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen  
527 Fan, Huayu Chen, Weize Chen, et al. The entropy mechanism of reinforcement learning for  
528 reasoning language models. *arXiv preprint arXiv:2505.22617*, 2025b.
- 529 Hanze Dong, Wei Xiong, Deep Goyal, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong  
530 Wang. Raft: Reward ranked finetuning for aligning language models with human feedback. *arXiv  
531 preprint arXiv:2304.06767*, 2023.
- 532 Kawin Ethayarajh, Lawrence Gao, and Dan Jurafsky. Kahneman-tversky optimization (kto): A new  
533 way to align language models. *arXiv preprint arXiv:2402.01306*, 2024.

- 540 Yuqian Fu, Tinghong Chen, Jiajun Chai, Xihuai Wang, Songjun Tu, Guojun Yin, Wei Lin, Qichao  
541 Zhang, Yuanheng Zhu, and Dongbin Zhao. Srft: A single-stage method with supervised and  
542 reinforcement fine-tuning for reasoning. *arXiv preprint arXiv:2506.19767*, 2025.
- 543 Amelia Glaese, Nat McAleese, Maja Mladenov, Sören Kaufmann, Amanda Askell, Phillip Butler,  
544 Tsim Chen, Courtney Voss, Vlad Cirroprocessing, Rachael Cummings, et al. Improving alignment of  
545 dialogue agents via targeted human judgements. *arXiv preprint arXiv:2209.14375*, 2022.
- 546 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad  
547 Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of  
548 models. *arXiv preprint arXiv:2407.21783*, 2024.
- 549 Caglar Gulcehre, Tom Jones, Ksenia Konyushkova, Florian Besse, David Budden, Angeliki Lazaridou,  
550 Son Nguyen, Razvan Dadashi, Jia He, et al. Reinforced self-training (rest) for language modeling.  
551 *arXiv preprint arXiv:2308.08998*, 2023.
- 552 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
553 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms  
554 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 555 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu,  
556 Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for  
557 promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint*  
558 *arXiv:2402.14008*, 2024.
- 559 Lixuan He, Jie Feng, and Yong Li. Amft: Aligning llm reasoners by meta-learning the optimal  
560 imitation-exploration balance. *arXiv preprint arXiv:2508.06944*, 2025a.
- 561 Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu Chen, Yue Wang, Linfeng Song, Dian  
562 Yu, Zhenwen Liang, Wenxuan Wang, et al. Deepmath-103k: A large-scale, challenging, de-  
563 contaminated, and verifiable mathematical dataset for advancing reasoning. *arXiv preprint*  
564 *arXiv:2504.11456*, 2025b.
- 565 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
566 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*  
567 *preprint arXiv:2103.03874*, 2021a.
- 568 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
569 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*,  
570 2021b.
- 571 Jian Hu, Jason Klein Liu, Haotian Xu, and Wei Shen. Reinforce++: An efficient rlhf algorithm with  
572 robustness to both prompt and reward models. *arXiv preprint arXiv:2501.03262*, 2025a.
- 573 Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xiangyu Zhang, and Heung-Yeung Shum.  
574 Open-reasoner-zero: An open source approach to scaling up reinforcement learning on the base  
575 model. *arXiv preprint arXiv:2503.24290*, 2025b.
- 576 Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec  
577 Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint*  
578 *arXiv:2412.16720*, 2024.
- 579 Sham Machandranath Kakade. *On the sample complexity of reinforcement learning*. University of  
580 London, University College London (United Kingdom), 2003.
- 581 Min-Joon Kim, Aviral Singh, and Hong-Seok Lee. Dynamic policy fusion for mixed-signal llm  
582 alignment. In *Third Conference on Language Modeling*, 2025.
- 583 Aviral Kumar, Joey Hong, Anikait Singh, and Sergey Levine. When should we prefer offline  
584 reinforcement learning over behavioral cloning? *arXiv preprint arXiv:2204.05618*, 2022.
- 585 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ra-  
586 masesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative  
587 reasoning problems with language models. *Advances in neural information processing systems*,  
588 35:3843–3857, 2022.

- 594 Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif  
595 Rasul, Longhui Yu, Albert Q Jiang, Ziju Shen, et al. NuminaMath: The largest public dataset in  
596 ai4maths with 860k pairs of competition math problems and solutions. *Hugging Face repository*,  
597 13:9, 2024.
- 598 Jia Li, Zhaofeng Wang, and Junxian He. Self-guided exploration with offline demonstrations for  
599 complex reasoning. In *Proceedings of the International Conference on Learning Representations*,  
600 2025.
- 601 Hao Liu, Zixuan Ji, and Di Lu. Bridging the gap between supervised fine-tuning and reinforcement  
602 learning. *arXiv preprint arXiv:2308.08809*, 2023.
- 603 Jason Liu, Zhiyuan Chen, and Ji-Woo Park. Direct fine-tuning on rewarded trajectories for language  
604 model alignment. *arXiv preprint arXiv:2406.13581*, 2024.
- 605 Jiazhen Liu, Yuchuan Deng, and Long Chen. Empowering small vlms to think with dynamic  
606 memorization and exploration. *arXiv preprint arXiv:2506.23061*, 2025a.
- 607 Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min  
608 Lin. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*,  
609 2025b.
- 610 Zihan Liu, Alekh Agarwal, and Nan Jiang. A principled analysis of offline preference optimization  
611 algorithms. *Journal of Machine Learning Research*, 26(45):1–58, 2025c.
- 612 Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V  
613 Le, Barret Zoph, Jason Wei, et al. The flan collection: Designing data and methods for effective  
614 instruction tuning. *arXiv preprint arXiv:2301.13688*, 2023.
- 615 Lu Ma, Hao Liang, Meiyi Qiang, Lexiang Tang, Xiaochen Ma, Zhen Hao Wong, Junbo Niu, Chengyu  
616 Shen, Runming He, Bin Cui, et al. Learning what reinforcement learning can’t: Interleaved online  
617 fine-tuning for hardest questions. *arXiv preprint arXiv:2506.07527*, 2025.
- 618 Panos Marantos, Yannis Koveos, and Kostas J Kyriakopoulos. Uav state estimation using adaptive  
619 complementary filters. *IEEE Transactions on Control Systems Technology*, 24(4):1214–1226,  
620 2015.
- 621 Eric Mitchell, Sergey Levine, and Chelsea Finn. Leveraging offline datasets for efficient online rl  
622 in large language models. In *Proceedings of the International Conference on Machine Learning*,  
623 2024.
- 624 Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher  
625 Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted  
626 question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- 627 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
628 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow  
629 instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:  
630 27730–27744, 2022.
- 631 Ji-Woo Park, Yifan Chen, and Denny Zhou. Reward-reweighted sft: An offline policy refinement  
632 method. *arXiv preprint arXiv:2502.11842*, 2025.
- 633 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Chelsea Finn, and Christopher D.  
634 Manning. Direct preference optimization: Your language model is secretly a reward model. *arXiv  
635 preprint arXiv:2305.18290*, 2023.
- 636 Neel Rajani, Aryo Pradipta, Gema Seraphina Goldfarb-Tarrant, and Ivan Titov. Scalpel vs. hammer:  
637 GRPO amplifies existing capabilities, SFT replaces them, 2025.
- 638 Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline rein-  
639 forcement learning and imitation learning: A tale of pessimism. *Advances in Neural Information  
640 Processing Systems*, 34:11702–11716, 2021.

- 648 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani,  
649 Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In  
650 *First Conference on Language Modeling*, 2024.
- 651
- 652 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region  
653 policy optimization. In *International conference on machine learning*, pp. 1889–1897. PMLR,  
654 2015a.
- 655 John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional  
656 continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*,  
657 2015b.
- 658
- 659 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
660 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 661
- 662 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,  
663 Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical  
664 reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 665
- 666 Aviral Singh, Joey Hong, and Aviral Kumar. Beyond reward: Offline preference-guided policy  
667 learning. In *Advances in Neural Information Processing Systems*, volume 36, 2023.
- 668
- 669 Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M Ziegler, Ryan Lowe, Christopher Hesse, John  
670 Schulman, and Jacob Hilton. Learning to summarize from human feedback. *Advances in Neural  
Information Processing Systems*, 33:3035–3045, 2020.
- 671
- 672 Wei Sun, Wen Yang, Pu Jian, Qianlong Du, Fuwei Cui, Shuo Ren, and Jiajun Zhang. Ktae: A model-  
673 free algorithm to key-tokens advantage estimation in mathematical reasoning. *arXiv preprint  
arXiv:2505.16826*, 2025.
- 674
- 675 Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun  
676 Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with  
677 llms. *arXiv preprint arXiv:2501.12599*, 2025.
- 678
- 679 Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. *arXiv preprint  
arXiv:1805.01954*, 2018.
- 680
- 681 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
682 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
683 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 684
- 685 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
686 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation  
687 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 688
- 689 Hugo Touvron, Louis Martin, and Guillaume Lample. Context distillation for on-policy reinforcement  
690 learning in llms. In *First Conference on Language Modeling*, 2024.
- 691
- 692 Shenzhi Wang, Le Yu, Chang Gao, Chujie Zheng, Shixuan Liu, Rui Lu, Kai Dang, Xionghui Chen,  
693 Jianxin Yang, Zhenru Zhang, et al. Beyond the 80/20 rule: High-entropy minority tokens drive  
694 effective reinforcement learning for llm reasoning. *arXiv preprint arXiv:2506.01939*, 2025.
- 695
- 696 Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Gu, Aitor Lewkowycz, Yao Lu, Ambrose Slone,  
697 Quoc Le, and Barret Zoph. Finetuned language models are zero-shot learners. *arXiv preprint  
arXiv:2109.01652*, 2021.
- 698
- 699 Jeff Wu, Long Ouyang, and Nisan Stiennon. Alternating between on-policy and off-policy updates  
700 for efficient and stable llm alignment. *arXiv preprint arXiv:2401.08543*, 2024.
- 701
- 702 Yongliang Wu, Yizhou Zhou, Zhou Ziheng, Yingzhe Peng, Xinyu Ye, Xinting Hu, Wenbo Zhu,  
Lu Qi, Ming-Hsuan Yang, and Xu Yang. On the generalization of sft: A reinforcement learning  
perspective with reward rectification. *arXiv preprint arXiv:2508.05629*, 2025.

- 702 Jianhao Yan, Yafu Li, Zican Hu, Zhi Wang, Ganqu Cui, Xiaoye Qu, Yu Cheng, and Yue Zhang.  
703 Learning to reason under off-policy guidance. *arXiv preprint arXiv:2504.14945*, 2025.  
704
- 705 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,  
706 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint*  
707 *arXiv:2407.10671*, 2024.
- 708 Zhicheng Yang, Zhijiang Guo, Yinya Huang, Xiaodan Liang, Yiwei Wang, and Jing Tang. Treerpo:  
709 Tree relative policy optimization. *arXiv preprint arXiv:2506.05183*, 2025.  
710
- 711 Hiroshi Yoshihara, Taiki Yamaguchi, and Yuichi Inoue. A practical two-stage recipe for mathematical  
712 llms: Maximizing accuracy with sft and efficiency with reinforcement learning. *arXiv preprint*  
713 *arXiv:2507.08267*, 2025.  
714
- 715 Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian  
716 Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system at  
717 scale. *arXiv preprint arXiv:2503.14476*, 2025.  
718
- 719 Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Shiji Song, and Gao Huang. Does  
720 reinforcement learning really incentivize reasoning capacity in llms beyond the base model? *arXiv*  
721 *preprint arXiv:2504.13837*, 2025.
- 722 Min Zeng, Jinfei Sun, Xueyou Luo, Caiquan Liu, Shiqi Zhang, Li Xie, and Xiaoxin Chen. Gta:  
723 Supervised-guided reinforcement learning for text classification with large language models. *arXiv*  
724 *preprint arXiv:2509.12108*, 2025a.  
725
- 726 Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. Simplerl-  
727 zoo: Investigating and taming zero reinforcement learning for open base models in the wild. In  
728 *Second Conference on Language Modeling*, 2025b.  
729
- 730 Weihao Zeng, Yuzhen Huang, Wei Liu, Keqing He, Qian Liu, Zejun Ma, and Junxian He. 7b model  
731 and 8k examples: Emerging reasoning with reinforcement learning is both effective and efficient.  
732 *HKUST-NLP Blog*, 2025c.
- 733 Wenhao Zhang, Yuexiang Xie, Yuchang Sun, Yanxi Chen, Guoyin Wang, Yaliang Li, Bolin Ding,  
734 and Jingren Zhou. On-policy rl meets off-policy experts: Harmonizing supervised fine-tuning and  
735 reinforcement learning via dynamic weighting. *arXiv preprint arXiv:2508.11408*, 2025.  
736
- 737 Weizhe Zhao, Benjamin Packer, and Ilya Kostrikov. Reward model fine-tuning using relative gradient  
738 updates. *arXiv preprint arXiv:2310.10574*, 2023.  
739
- 740 Xuandong Zhao, Zhewei Kang, Aosong Feng, Sergey Levine, and Dawn Song. Learning to reason  
741 without external rewards. *arXiv preprint arXiv:2505.19590*, 2025.
- 742 Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang,  
743 Yuqiong Liu, Rui Men, An Yang, et al. Group sequence policy optimization. *arXiv preprint*  
744 *arXiv:2507.18071*, 2025.  
745
- 746 Chunting Zhou, Graham Neubig, and Junxian He. Prefix-tuning for guided text generation in  
747 reinforcement learning. *Transactions of the Association for Computational Linguistics*, 11:1234–  
748 1249, 2023.  
749
- 750 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul  
751 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv*  
752 *preprint arXiv:1909.08593*, 2019.
- 753 Yuxin Zuo, Kaiyan Zhang, Li Sheng, Shang Qu, Ganqu Cui, Xuekai Zhu, Haozhan Li, Yuchen  
754 Zhang, Xinwei Long, Ermo Hua, et al. Ttrl: Test-time reinforcement learning. *arXiv preprint*  
755 *arXiv:2504.16084*, 2025.

## A MORE RELATED WORKS

### A.1 LLM POST-TRAINING: SFT AND RL

Current post-training methodologies for LLMs are largely centered around two primary paradigms: Supervised Fine-Tuning (SFT) and Reinforcement Learning (RL) (Wei et al., 2021; Ouyang et al., 2022). In the SFT paradigm, models are adapted for specific applications through training on curated input-output pairs, a process which has been shown to effectively align their behavior with human demonstrations (Chung et al., 2022; Longpre et al., 2023; Touvron et al., 2023a;b). In parallel, numerous works have highlighted RL as an effective approach for refining LLM behavior in ways that are difficult to capture with SFT’s static datasets (Glaese et al., 2022; Bai et al., 2022; Nakano et al., 2021). Within this domain, a popular framework is Reinforcement Learning from Human Feedback (RLHF), which optimizes the LLM policy against a reward model trained on human preferences (Christiano et al., 2017; Stiennon et al., 2020). Multiple works have established Proximal Policy Optimization (PPO) as a cornerstone algorithm for this phase (Schulman et al., 2017; Ziegler et al., 2019). To further improve reasoning capabilities in reward-driven optimization, recent advancements like Group Relative Policy Optimization (GRPO) have also been developed and widely adopted (Shao et al., 2024; Zheng et al., 2025; Chen et al., 2025).

### A.2 A COMBINATION OF ONLINE AND OFFLINE DATA IN LLM POST-TRAINING

Beyond applying SFT or RL in isolation, further explorations have sought to synergize their respective strengths by combining signals from pre-existing *offline data* and dynamically generated *online data* (Fu et al., 2025; Yan et al., 2025; Ma et al., 2025). This motivation stems from the distinct characteristics of each approach: SFT is noted for its efficiency in distilling knowledge from offline sources, whereas RL is valued for fostering exploration through online rollouts, a process frequently linked to improved generalization (Rajani et al., 2025; Chu et al., 2025). The strategies for this integration are diverse; some techniques use offline data as a prefix to guide online generation (Zhou et al., 2023; Touvron et al., 2024; Li et al., 2025; Wang et al., 2025), while others enhance offline data by incorporating reward signals in a process known as reward-augmented fine-tuning (Liu et al., 2024; Zhao et al., 2023; Park et al., 2025; Sun et al., 2025). The broader landscape also includes various purely offline preference optimization methods, though they follow a different paradigm (Rafailov et al., 2023; Mitchell et al., 2024; Liu et al., 2025c; Ethayarajh et al., 2024; Ahmadian et al., 2024; Wu et al., 2025). However, the most direct approach to synergy involves the concurrent use of both data types for training updates.

This direct approach, often termed mix-policy learning, is particularly relevant to our work and typically involves updating the model with a composite objective that combines an SFT loss from offline data and an RL loss from online data (Dong et al., 2023; Gulcehre et al., 2023; Singh et al., 2023; Liu et al., 2023). For instance, LUFFY (Yan et al., 2025) explores this paradigm by combining a fixed ratio of offline demonstration data with online rollouts in each training batch. Subsequently, SRFT (Fu et al., 2025) proposed a monolithic training phase that dynamically adjusts the weights of SFT and RL losses based on the model’s policy entropy, further demonstrating the viability of unifying these signals over a sequential pipeline. The principle of creating such a composite loss is shared by a variety of other recent frameworks (Wu et al., 2024; Zhang et al., 2025; Kim et al., 2025; Yu et al., 2025; Liu et al., 2025a; He et al., 2025a). While these methods highlight a clear trend towards unifying training signals, a foundational theoretical analysis explaining why these different learning signals can be effectively combined is still lacking. This motivates our work to establish a unified theoretical framework that in turn inspires a more principled algorithm design.

## B GRADIENT DERIVATION FOR CLASSICAL ALGORITHMS

### B.1 GRADIENT OF SFT

We first consider the SFT process as a warm-up. As mentioned in the previous section, SFT takes a pre-trained foundation model and further makes the model more specialized by training its output prediction distribution to align with domain-specific data. The fine-tuning process uses the same cross-entropy loss as in model pre-training, defined as follows,

810

811

812

813

$$\mathcal{L}_{SFT}(\theta) = - \sum_{i=1}^N \sum_{t=1}^{|\tau_i|} \log \pi_{\theta}(\tau_{i,t} | q_i, \tau_{i,<t}). \quad (12)$$

814

815

where  $\mathcal{D}_{SFT} = \{(q_i, \tau_i)\}_{i \in [N]}$  denotes the SFT dataset consisting of  $N$  question and trajectory pairs.  $\tau_t$  denotes the  $t$ -th token in the trajectory and  $\tau_{<t}$  denotes all the tokens prior to  $\tau_t$ .

816

817

818

819

For any  $t$ , the LLM outputs the next-token prediction as a probability distribution. In the context of RL, such a probability distribution has been commonly considered as a stochastic policy. Then, the gradient calculation of SFT can be obtained by directly taking the derivative of Equation (12) and takes the following form:

820

821

822

823

$$\nabla \mathcal{J}_{SFT}(\theta) = -\nabla \mathcal{L}_{SFT}(\theta) = \sum_{i=1}^N \sum_{t=1}^{|\tau_i|} \nabla \pi_{\theta}(\tau_{i,t} | q_i, \tau_{i,<t}) \frac{1}{\pi_{\theta}(\tau_{i,t} | q_i, \tau_{i,<t})}. \quad (13)$$

824

825

826

827

In this section, we slightly abuse the notion of policy gradient and consider the SFT as a case of behavioral cloning (BC) (Torabi et al., 2018), and Equation (13) can be seen as a specific form of policy gradient.

828

829

## B.2 GRADIENT OF ONLINE RL: PPO, GRPO AND BEYOND

830

831

832

833

834

835

836

For online RL, we first consider Proximal Policy Optimization (PPO) (Schulman et al., 2017) and a series of its derivations. PPO is a pivotal technique for RLVR in LLMs. Motivated by TRPO, PPO keeps the new policy close to the old policy, and performs conservative policy updates by incorporating a clipped version of its policy ratio in its objective. The clipping function was shown to stabilize the training process and avoid performance collapse during training. In this section, we omit the regularization terms, such as the KL divergence and entropy. The loss objective for PPO can be written as follows,

837

838

839

$$\mathcal{L}_{PPO}(\pi_{\theta}) = -\frac{1}{N} \sum_{i=1}^N \frac{1}{G} \sum_{j=1}^G \frac{1}{|\tau_j|} \sum_{t=1}^{|\tau_j|} \min(r_{i,j,t}(\theta) \hat{A}_{i,j}, \text{clip}(r_{i,j,t}(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_{i,j}), \quad (14)$$

840

841

842

843

844

845

846

847

848

849

In this setting, we consider questions sampled from a given dataset  $\mathcal{D}_{RL} \triangleq \{q_i\}_{i=1}^N$ , and for each question, we consider  $G$  trajectories independently sampled using a reference policy  $\pi_{ref}$ . We use  $r_{i,j,t}(\theta) = \frac{\pi_{\theta}(\tau_{i,j,t} | q_i, \tau_{i,j,<t})}{\pi_{ref}(\tau_{i,j,t} | q_i, \tau_{i,j,<t})}$  to denote the policy ratio  $\pi_{\theta} / \pi_{ref}$  introduced for importance sampling,  $\epsilon$  denotes the clipping factor for the importance sampling ratio, enhancing stability.

850

851

852

853

854

855

856

857

858

859

860

861

862

863

For PPO,  $\hat{A}$  is estimated using the Generalized Advantage Estimation (GAE) (Schulman et al., 2015b), calculated based on the reward of the sampled trajectories. For the case of GRPO, the advantage estimate  $\hat{A}$  is calculated based on a set of sampled trajectories. Given question  $q_i$ , a group of sampled roll-out trajectory  $\{\tau_{i,j}\}_{j \in [G]}$  with verifiable reward  $R(\tau_{i,j}) \in \{0, 1\}$ ,  $\hat{A}_{i,j}$  is calculated as the normalized reward over the group.

$$\hat{A}_{i,j} = \frac{R(\tau_{i,j}) - \text{mean}(\{R(\tau_{i,k})\}_{k \in [G]})}{\text{std}(\{R(\tau_{i,k})\}_{k \in [G]})}, \quad (15)$$

Compared to PPO, the most significant difference introduced by GRPO is the group relative advantage described above. Notably, the original manuscript of GRPO has also induced a sequence-level policy gradient balancing and a KL regularization term. However, more recent works such as (Yu et al., 2025) have removed or modified these terms in general.

The clipped surrogate objective in PPO and similar algorithms enhances the stability of the RL training process by turning off gradient propagation on samples where  $\pi_{\theta}$  moves too far from  $\pi_{ref}$ . For gradient calculation, this can be represented as an indicator function  $\mathbf{1}_{clip}$ .

$$\nabla \mathcal{J}_{PPO} = -\nabla \mathcal{L}_{PPO} = \frac{1}{N} \sum_{i=1}^N \frac{1}{G} \sum_{j=1}^G \frac{1}{|\tau_j|} \sum_{t=1}^{|\tau_j|} \nabla \pi_{\theta}(\tau_{i,j,t} | q_i, \tau_{i,j,<t}) \frac{\hat{A}_{i,j} \mathbf{1}_{clip}}{\pi_{ref}(\tau_{i,j,t} | q_i, \tau_{i,j,<t})}. \quad (16)$$

Apart from PPO and GRPO, many recent RL algorithms for RL post-training in LLMs can be shown to exhibit a similar form for their policy gradient calculations.

### B.3 GRADIENT OF OFFLINE RL

As stated in the previous sections, many recent studies seek to leverage offline data in the online RL training process for LLMs. These methods consider expert demonstration data as trajectories sampled from a near-optimal policy, and perform RL updates on these data based on policy gradient updates. These algorithms are adapted from the online RL literature and often combine offline and online training, setting them apart from simple SFT.

Taking SRFT (Fu et al., 2025) as an instance, the offline RL objective can be written as follows

$$\mathcal{L}_{SRFT}(\pi_\theta) = -\frac{1}{N} \sum_{i=1}^N \frac{1}{G} \sum_{j=1}^G \frac{1}{|\tau_j|} \sum_{t=1}^{|\tau_j|} \pi_\theta(\tau_{i,j,t} | q_i, \tau_{i,j,<t}) \hat{A}_{i,j}, \quad (17)$$

This objective is derived from the GRPO objective in Equation (14), while setting  $\pi_{ref} \equiv 1$  and removing the clipping mechanism since it becomes imbalanced. The motivation behind setting  $\pi_{ref} \equiv 1$  is that  $\pi_{ref}$  is typically unavailable for offline data. Under the assumption that the demonstration policy evenly covers the current policy  $\pi_\theta$ . In this case, setting  $\pi_{ref}$  to 1 changes the algorithm from importance sampling to rejection sampling. The policy gradient of the offline SRFT objective can be derived consequently.

$$\nabla \mathcal{J}_{SRFT} = -\nabla \mathcal{L}_{SRFT} = \frac{1}{N} \sum_{i=1}^N \frac{1}{G} \sum_{j=1}^G \frac{1}{|\tau_j|} \sum_{t=1}^{|\tau_j|} \nabla \pi_\theta(\tau_{i,j,t} | q_i, \tau_{i,j,<t}) \frac{\hat{A}_{i,j}}{\pi_{ref} = 1}. \quad (18)$$

## C GRADIENT COMPONENT ANALYSIS

### Takeaways

1. While all algorithms share the same Common Objective, bias-variance trade-offs still exist across current instances for different components of the unified gradient estimator.
2. We can improve the post-training process by constructing a better and more suitable estimation of the policy gradient.

Across the wide spectrum of algorithms contained in our previous discussions and Table 1, it can be inferred that the four components that construct the unified gradient estimator are motivated by different procedures in the post-training process. To better illustrate the relationship between the different processes with the respective components of our unified gradient, we present Figure 1.

We divide the post-training process of LLMs into the four steps shown in Figure 1: i) First, the LLM makes the decision on its data source, either to use data from an offline demonstration dataset, from self-generated rollout data, or a mixture of both. In this process, the policy likelihood  $\pi_\theta$  of the data with respect to the current LLM is generated. ii) Given the data source used for data generation, a reference policy  $\pi_{ref}$  is calculated. iii) After data collection is complete, the algorithm calculates the advantage estimation  $\hat{A}$  for each token/sequence. iv) Lastly, the algorithm may choose to apply an additional masking procedure  $\mathbb{1}_{stable}$  to disable the gradient calculation of various tokens, which could lead to theoretical or numerical stability issues. After these four steps, the components are collected to construct the policy gradient  $\text{grad}_{Uni}$ , which is used to update the LLM in the system. Similar to GAE presented in Schulman et al. (2015b), multiple instantiations exist to estimate the policy gradient. However, different component selections introduce various degrees of bias and variance, where a trade-off is often encountered. We provide the following discussion on key components of the unified gradient below.

**Reference Policy Calculation** Practically speaking, the reference policy denominator places a weight on each token-level update such that any token with a smaller probability, often implying more significance, is weighted more. SFT and REINFORCE assign weights inversely proportional

918 to the current policy  $\pi_\theta$ , enforcing a bigger update when the model outputs a small probability. On  
 919 the other hand, when the data is generated with an outdated model, algorithms such as PPO assign  
 920 weights inversely proportional to the rollout policy  $\pi_{\theta_{old}}$ , and offline RL does not assign additional  
 921 weights for tokens.

922 Theoretically, the reference policy is usually set given the source of the dataset and/or the rollout  
 923 policy. For online RL methods that train purely with on-policy data, such as REINFORCE (Ahmadian  
 924 et al., 2024), uses  $\frac{1}{\pi_\theta}$ , which produces an unbiased estimate for gradient calculation. However, these  
 925 methods usually suffer from high variance. For PPO-style online RL algorithms, the reference policy  
 926 refers to the rollout policy, which is a result of importance sampling. PPO is a numerically simplified  
 927 version of TRPO (Schulman et al., 2015a). PPO makes conservative updates that effectively reduce  
 928 variance. However, the important sampling ratio is in fact theoretically ill-posed and could introduce  
 929 systematic bias, as discussed in GSPO (Zheng et al., 2025). GSPO has also proposed a novel  
 930 calculation for  $\pi_{ref}$ , as shown in Table 1. On the other hand, in the offline setting, the choice for  
 931 reference policy  $\pi_{ref}$  is limited, since the algorithm generally has no access to the rollout policy. If  
 932 we are given the assumption that the offline data evenly covers the entire state-action rollout space,  
 933 then the importance sampling ratio  $r(\theta) = \frac{\pi_\theta(\tau)}{\pi_{ref}(\tau)}$  reduces to  $\pi_\theta(\tau)$  by setting constant  $\pi_{ref}(\tau) = 1$ .  
 934 Notably, it is apparent that setting  $\pi_{ref}(\tau) = 1$  introduces much bias at the cost of numerical stability.  
 935 For the SFT case, we can consider that the domain-specific dataset is generated with respect to the  
 936 expert policy  $\pi^*$ ; therefore, no weighted sampling is required. Neither of the two approaches is  
 937 entirely theoretically justified, from an RL perspective; both require a lower bound on the state-action  
 938 visitation of all the possible state-action pairs (Kakade, 2003), which can not be satisfied due to the  
 939 severely limited datasets in practice.

940 Apart from the strong connection to data source and sampling polices, some studies employ a hand-  
 941 crafted reweight factor within the reference policy denominator. These works (Yan et al., 2025; Zhang  
 942 et al., 2025) typically find desirable token properties and purposefully place a higher/lower weight on  
 943 these desirable/undesirable tokens, respectively.

944 **Choice of Stabilization Mask** The clipping operation introduced in PPO was the first to explicitly  
 945 add a stop gradient operation on LLM post-training. Clipping gradient estimation where the impor-  
 946 tance sampling strays too far from 1 is an effective approach to address high variances. However,  
 947 this aggressive clipping behavior has been criticized by some to be overly conservative: Both DAPO  
 948 (Yu et al., 2025) and CISPO (Chen et al., 2025) stated that the classical PPO approach drops all  
 949 the tokens corresponding to large model updates, and that many such tokens are in fact crucial for  
 950 stabilizing entropy and facilitating scalable RL. DAPO presented a slight modification to the clipping  
 951 threshold, and CISPO further extended the notion of token-wise mask, where more granular tuning  
 952 was introduced to decide whether gradients from specific tokens should be dropped. The recent  
 953 work of Cui et al. (2025b) has demonstrated that many existing algorithms negatively impact the  
 954 output entropy during training and introduced Clip-Cov, adding another clipping mechanism to  
 955 address the entropy-collapse encountered in training. While these methods demonstrated performance  
 956 enhancements in practice, they also provide additional sources of bias.

957 On the other hand, works such as GSPO (Zheng et al., 2025) have stated that the PPO-style clipping  
 958 is inherently noisy and inefficient for sample exploitation: GSPO clips a much larger fraction of  
 959 tokens and yet demonstrated superior training efficiency.

960 In addition, post-training algorithms using offline data have chosen to purposefully remove the  
 961 clipping from training, mostly guided by performance. Though setting  $\pi_{ref}(\tau) = 1$  as the policy  
 962 denominator does effectively reduce the instability in gradient calculations.

964 **Advantage Estimation** There are two commonly used settings for estimating the sequence-level  
 965 advantage function: the fixed advantage setting and the adaptive advantage setting. The fixed setting  
 966 considers  $\hat{A} = \pm 1$  given the rule-based verification, which is adapted by REINFORCE and implicitly  
 967 by SFT (where all sequences are positive samples). Alternatively, recent studies have focused on using  
 968 adaptive advantage estimations, performing re-centering or normalization based on the performance  
 969 of the current rollout group. Notably, GRPO and its variants, such as DAPO (Yu et al., 2025) and  
 970 LUFFY (Yan et al., 2025), use unit normalization such that the advantage estimation of the group  
 971 has a unit standard deviation. Other approaches, such as Dr. GRPO (Liu et al., 2025b), RLOO  
 (Ahmadian et al., 2024), and REINFORCE++ (Hu et al., 2025a), claim that dividing the standard

972 deviation introduces a difficulty bias and that only recentering is adequate.

973  
974 Apart from sequence-level advantage estimate  $\hat{A}_{i,j}$ , recent works (Wang et al., 2025; Yang et al.,  
975 2025; Sun et al., 2025) have also adapted a more granular token-level advantage estimate  $\hat{A}_{i,j,t}$  to a  
976 varying degree of success.

977  
978 **A Combination of Gradient Estimators** Although bias-variance trade-offs exist for the gradient  
979 estimator, we state that, given data distribution assumptions and sufficient data samples, all policy  
980 gradient estimators covered in our framework should result in an effective direction of improvement  
981 for the Common Objective. To effectively reduce the variance and bias for each policy update, we  
982 can treat instances of policy gradient as different noisy measurements of the true policy gradient, and  
983 perform a weighted average to generate a more accurate gradient estimation, similar to complementary  
984 filters (Marantos et al., 2015).

985 However, the complexity of LLM RLVR introduces additional challenges. The current state of the  
986 behavior policy  $\pi_\theta$  and its relationship with the respective tasks also greatly impacts the bias-variance  
987 tradeoff of each instance of the gradient estimator. For instance, RL-zero is significantly more  
988 effective for the Qwen model series compared to LLaMA, but SFT is effective for both methods (Zeng  
989 et al., 2025b); SFT  $\rightarrow$  RL and RL  $\rightarrow$  SFT also yield significantly different results on the same LLM  
990 (Fu et al., 2025). We argue that for constructing a post-training algorithm with better effectiveness and  
991 efficiency, a dynamic and adaptive mechanism is crucial to construct optimal gradient components.

## 992 D ADDITIONAL THEORETICAL DETAILS FOR SECTION 2.2

### 993 D.1 DERIVING EQUATION 2 FROM EQUATION 1

994  
995 **Lemma A1** (Score-function identity). *For density  $\pi_\theta$  and integrable  $f(\tau)$ ,*

$$996 \quad \nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta} [f(\tau)] = \mathbb{E}_{\tau \sim \pi_\theta} [f(\tau) \nabla_\theta \log \pi_\theta(\tau)], \quad \mathbb{E}_{\tau \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(\tau)] = 0.$$

997  
998 **Lemma A2** (Differentiating an expectation with parameterized integrand). *For differentiable  $f_\theta$ ,*

$$999 \quad \nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta} [f_\theta(\tau)] = \mathbb{E}_{\tau \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(\tau) f_\theta(\tau) + \nabla_\theta f_\theta(\tau)].$$

1000  
1001 **Lemma A3** (Measure-change (importance reweighting) identity). *Let  $s(\tau | q)$  be any sampling  
1002 density that is positive wherever  $\pi_\theta(\tau | q)$  is. Then*

$$1003 \quad \mathbb{E}_{\tau \sim \pi_\theta} [f(\tau) \nabla_\theta \log \pi_\theta(\tau)] = \mathbb{E}_{\tau \sim s} \left[ \frac{\pi_\theta(\tau)}{s(\tau)} f(\tau) \nabla_\theta \log \pi_\theta(\tau) \right] = \mathbb{E}_{\tau \sim s} \left[ \frac{1}{s(\tau)} f(\tau) \nabla_\theta \pi_\theta(\tau) \right].$$

1004  
1005 *Proof.* By Lemma A1,  $\nabla \mathbb{E}_{\pi_\theta} [r(\cdot | q)] = \mathbb{E}_{\pi_\theta} [r(\cdot | q) \nabla \log \pi_\theta]$ . For the data-adherence term, since  
1006  $\text{KL}(\pi_\beta | \pi_\theta) = \mathbb{E}_{\pi_\beta} [\log \pi_\beta - \log \pi_\theta]$  and  $\pi_\beta$  does not depend on  $\theta$ , we have  $-\mu \nabla \text{KL}(\pi_\beta | \pi_\theta) =$   
1007  $\mu \mathbb{E}_{\pi_\beta} [\nabla \log \pi_\theta]$ . Summing yields the claim.  $\square$

### 1008 D.2 EXTENSION: ADDING A TRUST-REGION REGULARIZER

1009  
1010 A trust region encourages conservative policy updates by penalizing the KL divergence from the  
1011 current policy  $\pi_\theta$  to a fixed reference policy  $\pi_{ref}$ :

$$1012 \quad \lambda \text{KL}(\pi_\theta(\cdot | q) \| \pi_{ref}(\cdot | q)), \quad \lambda \geq 0.$$

1013  
1014 It is the penalty form of the constrained problem

$$1015 \quad \max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta} [r(\tau | q)] \quad \text{s.t.} \quad \text{KL}(\pi_\theta \| \pi_{ref}) \leq \delta,$$

1016  
1017 where  $\lambda$  acts as the Lagrange multiplier tied to the trust-region radius  $\delta$ . Typical choices are  
1018  $\pi_{ref} = \pi_{\theta_{old}}$  (on-policy stability, TRPO/PPO-style). This penalty controls step sizes, dampens  
1019 distribution shift, and yields clipping-style masks when optimized with PPO surrogates.

1020  
1021 **Objective and gradient with trust region.** Augmenting the Common Objective with the trust-region  
1022 term gives

$$1023 \quad \tilde{\mathcal{J}}_{\lambda, \mu}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta(\cdot | q)} [r(\tau | q)] - \lambda \text{KL}(\pi_\theta(\cdot | q) \| \pi_{ref}(\cdot | q)) - \mu \text{KL}(\pi_\beta(\cdot | q) \| \pi_\theta(\cdot | q)),$$

whose gradient is

$$\nabla_{\theta} \tilde{\mathcal{J}}_{\lambda, \mu}(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ (r(\tau | q) - \lambda \log \frac{\pi_{\theta}(\tau | q)}{\pi_{ref}(\tau | q)}) \nabla_{\theta} \log \pi_{\theta}(\tau | q) \right] + \mu \mathbb{E}_{\tau \sim \pi_{\beta}} [\nabla_{\theta} \log \pi_{\theta}(\tau | q)].$$

In the estimator (Equation 3), this corresponds to replacing the unified advantage by

$$\hat{A}_{uni}^{(\lambda)}(\tau, q) = r(\tau | q) - \lambda \log \frac{\pi_{\theta}(\tau | q)}{\pi_{ref}(\tau | q)} + \mu \mathbb{1}\{\pi_{ref} = \pi_{\beta}\} \frac{\pi_{\beta}(\tau | q)}{\pi_{\theta}(\tau | q)}.$$

All other expressions, including the masked estimator in Equation 5, remain unchanged in form (with  $\hat{A}_{uni}$  replaced by  $\hat{A}_{uni}^{(\lambda)}$ ).

### D.3 PPO CLIPPING AND THE STABILIZATION MASK

With rollout policy  $\pi_{\theta_{old}}$  and trust-region constraint  $\text{KL}(\pi_{\theta} \parallel \pi_{\theta_{old}}) \leq \delta$ , the PPO surrogate

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta_{old}}} \left[ \min (r_{\theta}(\tau) A_{\theta_{old}}(\tau), \text{clip}(r_{\theta}(\tau), 1 - \epsilon, 1 + \epsilon) A_{\theta_{old}}(\tau)) \right], \quad r_{\theta}(\tau) = \frac{\pi_{\theta}(\tau)}{\pi_{\theta_{old}}(\tau)},$$

has a piecewise derivative that is zero outside the trusted region in the harmful direction, yielding

$$\nabla_{\theta} \approx \mathbb{E}_{\tau \sim \pi_{\theta_{old}}} \left[ \mathbb{1}_{stable}(\tau) \frac{1}{\pi_{\theta_{old}}(\tau)} A_{\theta_{old}}(\tau) \nabla_{\theta} \pi_{\theta}(\tau) \right], \quad (19)$$

which matches the masked Unified Policy Gradient Estimator with  $\pi_{ref} = \pi_{\theta_{old}}$  and  $\hat{A} = A_{\theta_{old}}$ .

## E THE PROCEDURE AND ANALYSIS OF HYBRID POST-TRAINING (HPT)

---

### Algorithm 1 The Hybrid Post-Training (HPT) Algorithm

---

**Input:** Pretrained LLM (policy)  $\pi_{\theta}$ ; SFT dataset  $\mathcal{D}_{\text{SFT}} = \{(q, \tau^*)\}$  with supervising trajectories  $\tau^*$ ; verifier  $v$ ; on-policy samples number  $n$ ; total training steps  $T$ ; feedback functions  $f$  and  $g$ ; learning rate  $\eta$

**Output:** Fine-tuned policy  $\pi_{\theta^*}$ .

**for**  $t = 1$  **to**  $T$  **do**

**for**  $i = 1$  **to**  $n$  **do**

    Sample trajectory  $\tau_i \sim \pi_{\theta}(\cdot | q)$  Evaluate with verifier (rule-based reward):  $v(\tau_i) \leftarrow R(\tau_i) \in \{0, 1\}$

**end**

$P \leftarrow \frac{1}{n} \sum_{i=1}^n v(\tau_i)$   $\alpha \leftarrow f(P)$ ,  $\beta \leftarrow g(P)$  # Performance feedback on  $q$

  Compute RL loss  $\mathcal{L}_{\text{RL}}$  using rollouts  $\{\tau_i\}$  and normalized advantages derived from  $\{R(\tau_i)\}$ .

  Compute SFT loss  $\mathcal{L}_{\text{SFT}}$  on the supervising trajectory  $\tau^*$ .

$\mathcal{L} \leftarrow \alpha \mathcal{L}_{\text{RL}} + \beta \mathcal{L}_{\text{SFT}}$  # Mixed loss with feedback coefficients

$\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}$

**end**

**return**  $\pi_{\theta^*}$

---

We provide a detailed analysis of the effectiveness and advantages of HPT, the procedure of which is shown above.

**HPT as an Optimal Gradient Selection Mechanism** As stated in Section 2, we can consider the gradient calculations of existing approaches (including SFT, offline RL, and online RL) as different instances of gradient calculation, shown in Table 1. We argue that, similar to the choice between behavior cloning and reinforcement learning in robotics learning (Kumar et al., 2022), the exact gradient calculation should be carefully chosen for the post-training process of reasoning tasks.

When the model’s current policy is far away from the optimal policy, SFT, or behavior cloning, should be preferred over RL. Usually, this is indicated by the poor performance of the language model. This is evidently shown in the Reference Policy of SFT in Table 1. The current model policy being far from the demonstration would result in the Reference Policy  $\pi_{\theta}(\tau) \ll 1$ . Therefore, the SFT

1080 gradient could make a significantly larger model update compared to RL methods. Utilizing SFT  
 1081 could speed up the post-training process and increase performance.

1082 On the other hand, when the model exhibits relatively good performance, it implies that the model’s  
 1083 current policy is close to the optimal policy, and the online state-action distribution could more  
 1084 effectively cover the optimal response. Theoretically, this can be represented by the concentrability  
 1085 coefficient (Rashidinejad et al., 2021). In this case, online RL is preferred over SFT for its scalability.  
 1086 In addition, (Kumar et al., 2022) has shown that for problems with long-horizon properties and the  
 1087 environment has “critical” states. Both properties are satisfied in our math reasoning task, given  
 1088 recent studies such as Wang et al. (2025).  
 1089

1090 **Amplification of Success Probability via HPT** By reverting to SFT when the model fails entirely,  
 1091 we ensure policy improvement by injecting the high-quality external signal. When the model already  
 1092 achieves partial success, RL-style updates amplify the probability of success, which is similar to the  
 1093 success amplification property analyzed in RL theory. Therefore, for initial policy success rate  $p_0$ ,  
 1094 HPT induces iterative improvement that satisfies:

$$p_{t+1} \geq p_t + \delta(p_t),$$

1095 where in failure cases, supervised injection immediately boosts  $p_t$ , while in success cases the RL  
 1096 update yields monotonic amplification. By combining both, HPT retains a guaranteed upward drift in  
 1097 success probability.  
 1098  
 1099

## 1101 F DETAILS OF THE EXPERIMENTAL SETUP

1102  
 1103 For the benchmarks, AMC (Li et al., 2024) comprises problems drawn from the AMC12 2022 and  
 1104 AMC12 2023 examinations. ARC-c (Clark et al., 2018) is an open-domain reasoning benchmark and  
 1105 GPQA-Diamond (Rein et al., 2024) is a challenging and high-quality subset of the Graduate-Level  
 1106 Google-Proof Question Answering benchmark.

1107 We set the training batch size to 128 and the maximum generation length to 8, 192 tokens, unless  
 1108 otherwise specified. We keep the system prompt and training datasets the same as LUFFY (Yan et al.,  
 1109 2025) and perform evaluation on DeepMath (He et al., 2025b). For other details that may not have  
 1110 been explicitly introduced, we have endeavored to follow previous works as closely as possible (Zhao  
 1111 et al., 2025; Zuo et al., 2025). All experiments were conducted on 8 x NVIDIA A800 80GB GPUs.  
 1112

## 1113 G ABLATION STUDIES

1114  
 1115 For the ablation studies, § G.1 explores the role of off-policy RL, testing whether alternative strategies  
 1116 for utilizing offline data yield benefits. § G.2 presents a gate threshold ablation study.  
 1117

### 1118 G.1 IMPACT OF OFF-POLICY RL

1119  
 1120 Table 5: Performance of different training paradigms to evaluate the impact of Off-policy RL. **SFT/ON** denotes  
 1121 SFT/On-policy (HPT), **OFF/ON** denotes Off-policy/On-policy, and **Mix/ON** denotes Mix-policy/On-policy.

1122 Name	AIME 24	AIME 25	AMC	MATH-500	Minerva	Olympiad	Avg
1123 OFF/ON	16.6	11.8	47.3	76.2	35.3	41.6	38.1
1124 Mix/ON	<b>16.7</b>	17.2	46.9	79.4	37.5	43.9	40.3
1125 SFT/ON	16.6	<b>17.8</b>	<b>51.0</b>	<b>81.0</b>	<b>37.5</b>	<b>47.3</b>	<b>41.9</b>

1126  
 1127 In our previous work, we have only made preliminary attempts at unifying post-training by integrating  
 1128 RL with SFT. However, off-policy RL represents an important training paradigm that emphasizes  
 1129 leveraging the offline data. To this end, we further conduct experiments to investigate its influence  
 1130 and potential role.  
 1131

1132 We compare three different training paradigms: (1) SFT/On-policy, the model alternates between  
 1133 SFT and on-policy RL, which corresponds to the method we introduced above (HPT); (2) Off-  
 policy/On-policy, the model alternates between off-policy RL and on-policy RL during training; and

(3) Mix-policy/On-policy, the model combines the loss from SFT and off-policy RL, and dynamically switches it with the on-policy RL objective. For the Mix setting, we performed hyperparameter search and found the optimal SFT/OFF weighting ratio to be 1/10, i.e., the coefficients of the SFT loss and the off-policy loss are set to 0.1 and 1.0, respectively. We replicate the off-policy RL implementation described in LUFFY (Yan et al., 2025), and all experiments are conducted in the same settings to ensure fairness.

We evaluate the results of three methods on six math benchmarks. Table 5 presents results. Overall, the SFT/ON method achieves the best average performance (41.9), outperforming both Mix/ON (40.3) and OFF/ON (38.1). This suggests that off-policy RL may not be essential, as SFT already serves effectively as the training method of HPT for learning from offline data.

## G.2 GATE THRESHOLD ABLATION

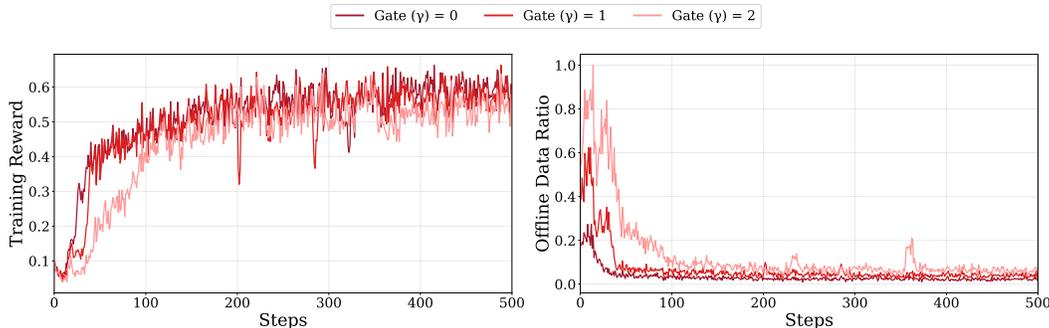


Figure 8: Training reward (left) and offline data ratio (right) comparisons across different gate settings on Qwen2.5-Math-1.5B.

In this section, we investigate the effect of different gate thresholds  $\gamma$ . A value of  $\gamma = 0$  indicates that the model switches to SFT only when it fails all questions. Similarly,  $\gamma = 1$  and  $\gamma = 2$  correspond to settings where the model remains in on-policy reinforcement learning as long as it answers at least one or two out of eight questions correctly, respectively. To visualize the impact of the gating mechanism, we conduct experiments on the Qwen2.5-Math-1.5B under three different gate settings. As shown in Figure 8, we analyze the training dynamics by tracking the dynamics of rewards and the proportion of offline data utilized throughout training, thereby highlighting how different gate thresholds mediate the balance between leveraging offline demonstrations and incorporating online feedback. We observe that, under different gate thresholds, varying degrees of engagement with offline data-based SFT learning emerge. A larger gate threshold introduces a greater extent of SFT based on offline data, as expected.

Table 6: Performance of HPT with different switch gate  $\gamma$  on Qwen2.5-Math-1.5B.

Name	AIME 24	AIME 25	AMC	MATH-500	Minerva	Olympiad	Avg
$\gamma = 2$	15.8	13.0	49.0	77.6	34.6	44.1	39.0
$\gamma = 1$	<b>18.1</b>	14.2	46.0	75.4	35.7	42.5	38.7
$\gamma = 0$	16.6	<b>17.8</b>	<b>51.0</b>	<b>81.0</b>	<b>37.5</b>	<b>47.3</b>	<b>41.9</b>

To further compare the performance across different gating strategies, we evaluate the three trained models on six benchmarks. Table 6 presents the results. Among the three configurations,  $\gamma = 0$  achieves the best overall performance with an average score of 41.9, outperforming both  $\gamma = 1$  (38.7) and  $\gamma = 2$  (39.0). This observation suggests that simply incorporating more SFT does not necessarily lead to better outcomes. Instead, it is crucial to maintain a dynamic balance between the *exploration* of RL and the *exploitation* of SFT. The optimal degree of this gating mechanism should be adjusted according to the characteristics of the base model and the specific training data employed.

1188 H THE USE OF LARGE LANGUAGE MODELS

1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

We use large language models to polish our writing. Specifically, we employ ChatGPT (GPT-5 Thinking) to revise our initial manuscript section by section. The prompt we use is *I am writing an academic paper in English. Please polish the following draft so that it adheres to the conventions of academic writing.*