

PROXIMAL SUPERVISED FINE-TUNING

Wenhong Zhu^{1,2} Ruobing Xie^{3,*} Rui Wang^{1,2} Xingwu Sun^{3,4} Di Wang³ Pengfei Liu^{1,2,*}

¹Shanghai Jiao Tong University ²Shanghai Innovation Institute ³Tencent

⁴University of Macau

{zwhong714, wangrui12, pengfei}@sjtu.edu.cn

{xrbsnowing}@163.com

ABSTRACT

Supervised fine-tuning (SFT) of foundation models often leads to poor generalization, where prior capabilities deteriorate after tuning on new tasks or domains. Inspired by trust-region policy optimization (TRPO) and proximal policy optimization (PPO) in reinforcement learning (RL), we propose **Proximal SFT (PSFT)**, a fine-tuning objective that incorporates the benefits of trust-region, effectively constraining policy drift during SFT while maintaining competitive tuning. By viewing SFT as a special case of policy gradient methods with constant positive advantages, we derive PSFT that **stabilizes optimization and leads to generalization**, while **leaving room for further optimization in subsequent post-training stages**. Experiments across mathematical and human-value domains show that PSFT matches SFT in-domain, outperforms it in out-of-domain generalization, remains stable under prolonged training without causing entropy collapse, and provides a stronger foundation for the subsequent optimization.¹

1 INTRODUCTION

Recently, post-training has become a crucial part of the overall training process. In particular, reinforcement learning (RL) algorithms, such as PPO (Schulman et al., 2017) and GRPO (Shao et al., 2024), have demonstrated significant effectiveness when applied to language models (LMs) focused on reasoning tasks. As RL is scaled over time, foundation models gain the capacity to address complex problems through more profound and extended reasoning (OpenAI, 2024; Guo et al., 2025). These reasoning models offer an abundant and valuable latent thoughts (Ruan et al., 2025) across the internet. Numerous community efforts have focused on leveraging this knowledge via supervised fine-tuning (SFT) (Guha et al., 2025; Li et al., 2025), a distillation approach valued for its efficiency and simplicity compared to RL.

However, SFT models are often criticized for poor generalization (Huan et al., 2025). This limitation arises because SFT essentially performs behavior cloning, which can result in weak generalization when the fine-tuning dataset is suboptimal or distributionally misaligned with the pretraining data (Chu et al., 2025), potentially causing large policy updates (Schulman et al., 2015). RL fine-tuning (RFT) provides a promising alternative, as numerous studies have shown that RL better preserves the generalization capabilities that SFT tends to erode (Chu et al., 2025; Huan et al., 2025).

Another concern is maintaining the ability to explore. In practice, SFT is often used as a cold start to stabilize RFT training. The gains from RFT may largely come from refining capabilities already acquired during pretraining and SFT (Gandhi et al., 2025). However, excessive reliance on SFT can diminish a model’s capacity for exploration (Xie et al., 2024), as it would cause the entropy collapse (Cui et al., 2025) and thereby constrain exploration during RL training (Yu et al., 2025).

Therefore, the current challenges lie in improving the **generalization** and **exploration keeping** of SFT models. To tackle these challenges, this paper proposes an improved fine-tuning strategy **Proximal Supervised Fine-Tuning (PSFT)** that avoids rote learning and achieves reliable performance. Specifically, we establish a theoretical connection between SFT and RL, and

*Corresponding authors.

¹Code is available at <https://github.com/zwhong714/PSFT>.

introduce a novel objective based on a clipped surrogate objective that leverages the benefits of trust regions to constrain policy updates.

Our experiments demonstrate that, compared to standard SFT, PSFT attains comparable performance on the target task while preserving the model’s general capabilities. Moreover, PSFT mitigates entropy collapse during training and yields superior target and generalization performance in subsequent RL stages. We evaluate PSFT on both the mathematics and human value alignment domains, highlighting the broad applicability of PSFT in practice.

Our contributions are as follows:

- We propose PSFT, an optimization method for SFT that adopts a clipped surrogate objective, akin to PPO, to enforce trust-region-like constraints and thereby prevent excessive policy update. This approach maximally preserves the model’s general capabilities while maintaining comparable performance on the target task.
- PSFT effectively prevents entropy collapse and overfitting during SFT, thereby enabling subsequent RL stages to achieve more robust and superior results on both target-specific and general tasks.
- We extensively validate PSFT across various base models, tasks, and evaluation metrics, demonstrating its potential as a promising alternative to standard SFT.

2 PRELIMINARIES AND MOTIVATIONS

2.1 MARKOV DECISION PROCESS FORMULATION

We formulate the problem of language modeling and optimization in the framework of a Markov Decision Process (MDP). An MDP is defined by a tuple $(\mathcal{S}, \mathcal{A}, P)$, where: \mathcal{S} is the state space, representing partial sequences. \mathcal{A} is the action space, representing possible next tokens. $P(s' | s, a)$ is the transition probability of moving from state s to s' given action a . In the context of autoregressive LMs, a query sequence $x := (x_1, \dots, x_m) \in \mathcal{X}$ initializes the state. A response sequence $y := (y_1, \dots, y_n) \in \mathcal{Y}$ is generated by sequentially sampling actions (tokens) from the policy π_θ . The joint probability of y given x is: $\pi_\theta(y | x) = \prod_{t=1}^n \pi_\theta(y_t | y_{<t}, x)$, where $y_{<t} := \{y_1, \dots, y_{t-1}\}$ and the state at time t is $s_t = (x, y_{<t})$.

2.2 SFT AS THE SPECIAL POLICY UPDATE

Supervised Fine-tuning. The training objective is to minimize the cross-entropy loss between the model’s predicted token distribution and the ground truth tokens. Formally, the loss is defined as:

$$L^{\text{SFT}}(\theta) = -\hat{\mathbb{E}}_{(s_t, a_t^*) \sim \mathcal{D}} [\log \pi_\theta(a_t^* | s_t)], \quad (1)$$

where (s_t, a_t^*) pairs are sampled from an offline dataset \mathcal{D} .

Policy Gradient. In contrast, policy gradient directly samples trajectories from the current policy π_θ interacting with the environment. Using the policy gradient theorem, the corresponding loss function can be defined as follows:

$$L^{\text{PG}}(\theta) = \hat{\mathbb{E}}_{(s_t, a_t) \sim \pi_\theta} [\log \pi_\theta(a_t | s_t) \hat{A}_t], \quad (2)$$

where \hat{A}_t is the estimated advantage function at step t .

From this perspective, SFT can be seen as a special case of policy gradient where sampling is from a fixed offline dataset \mathcal{D} and the advantage is fixed as $\hat{A}_t = 1$ for the ground-truth action, which corresponds to maximizing the likelihood of expert tokens.

2.3 PROXIMAL POLICY OPTIMIZATION

Trust Region Policy Optimization. TRPO (Schulman et al., 2015) maximizes a surrogate objective by leveraging trajectories collected under the previous policy $\pi_{\theta_{\text{old}}}$, introducing importance sampling to reweight these samples for evaluating the new policy π_θ :

$$L^{\text{CPI}}(\theta) = \hat{\mathbb{E}}_{(s_t, a_t) \sim \pi_{\theta_{\text{old}}}} \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right] = \hat{\mathbb{E}}_t \left[r_t(\theta) \hat{A}_t \right], \quad (3)$$

where $r_t(\theta)$ is the importance sampling ratio, and the superscript *CPI* refers to conservative policy iteration (Kakade & Langford, 2002). By introducing a trust region constraint on the KL divergence between π_{θ} and $\pi_{\theta_{\text{old}}}$, TRPO ensures that each policy update does not deviate too far from the reference policy.

Proximal Policy Optimization. However, directly maximizing L^{CPI} with a hard KL constraint, as in TRPO, can be difficult to optimize in practice. To address this, PPO (Schulman et al., 2017) modifies the surrogate objective by penalizing policy changes that move $r_t(\theta)$ too far from 1. The clipped surrogate objective is:

$$L^{\text{CLIP}}(\theta) = \hat{\mathbb{E}}_{(s_t, a_t) \sim \pi_{\theta_{\text{old}}}} \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right]. \quad (4)$$

This clipping mechanism effectively defines a soft trust region: $r_t(\theta)$ is restricted to remain within a small neighborhood of 1, which prevents destructive updates. Compared to TRPO, PPO achieves similar stability and efficiency.

3 PROXIMAL SUPERVISED FINE-TUNING

Proximal Supervised Fine-Tuning aims to improve performance on the target SFT task while preserving general capabilities, preventing entropy collapse, and allowing further optimization.

3.1 FROM PPO TO PSFT

We revisit the TRPO and PPO objectives in RL. Both methods rely on importance sampling ratios $r_t(\theta)$ to reweight returns, while constraining these ratios to prevent the new policy π_{θ} from deviating excessively from the old policy $\pi_{\theta_{\text{old}}}$. Translating this idea to the supervised setting, where all actions are assumed to be “correct” (i.e., $\hat{A}_t > 0$), we simplify the advantage to $\hat{A}_t = 1$, and define the Proximal SFT loss as:

$$L^{\text{PSFT}}(\theta) = \mathbb{E}_{(s_t, a_t) \sim \mathcal{D}} \left[\min \left(\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}, \text{clip}\left(\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}, 1 - \epsilon, 1 + \epsilon\right) \right) \right]. \quad (5)$$

This objective regularizes supervised learning updates by limiting the ratio between the new and old policy probabilities. We allow the old policy $\pi_{\theta_{\text{old}}}$ —rather than the fixed initial model—to evolve dynamically. Intuitively, it discourages overconfident changes in token probabilities, thereby preserving existing capabilities. When changing to the online setting, this method is equivalent to standard SFT training.

Warm-Up. Note that (s_t, a_t) in Eq. 5 is sampled from the offline dataset \mathcal{D} . In the initial steps, $r_t(\theta)$ may yield a biased expectation since $\pi_{\theta_{\text{old}}}$ is not aligned with the distribution of \mathcal{D} . The clipping mechanism constrains policy updates within an acceptable range, which may indirectly reduce the negative impact of this mismatch. A warm-up SFT phase on \mathcal{D} can be introduced to better align the initial policy $\pi_{\theta_{\text{old}}}$ with the offline dataset, further improving the SFT/PSFT in-domain performance.

3.2 GRADIENT ANALYSIS

The gradient updated is confined to the trust region to perform fine-tuning, similar to that in PPO. The gradient is as follows:

$$\begin{aligned} \nabla_{\theta} L^{\text{PSFT}}(\theta) &= \mathbb{E}_{(s_t, a_t) \sim \mathcal{D}} \left[r_t \cdot \mathbb{I}_{\text{trust}}(r_t) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right], \\ \text{where } \mathbb{I}_{\text{trust}}(r_t) &= \begin{cases} 0 & r_t > 1 + \epsilon, \\ 1 & \text{otherwise.} \end{cases} \end{aligned} \quad (6)$$

The gradient analysis shows that if the train dataset distribution deviates significantly from the model distribution, these tokens will have no gradient, thereby avoiding large policy updates and maintaining generalization. Typically, the optimal value of ϵ is set to 0.2 or 0.28. A larger ϵ can lead to large gradients, as shown in Equation 6. Further analysis is provided in Section 5.2.

4 EXPERIMENTS

In Section 4.1, we explore the training dynamics of our PSFT with standard SFT, and then evaluate both in-domain and out-of-domain performances of models fine-tuned with each method, focusing primarily on the mathematics domain. In Section 4.2, we attempt to verify the effectiveness of PSFT as the cold start point of RL. In Section 4.3, we further conduct experiments in the alignment domain to assess the universal applicability of our approach.

4.1 MAIN EXPERIMENTS ON MATH REASONING IN THE SFT STAGE

Setup. (1) *Models and Datasets.* We evaluate our method based on Qwen2.5-7B-Instruct (Yang et al., 2025) and Llama3.1-8B-Instruct (Dubey et al., 2024). Our training data focuses on the math domain, aiming to improve general LLM capabilities through math reasoning. We employ the OpenR1-Math-8192² long chain-of-thought (CoT) dataset (Face, 2025). (2) *Baseline.* We consider the SFT method and an SFT variant, denoted as SFT_{KL}, which incorporates KL divergence constraints into the loss function. The KL regularization coefficient is set to 0.5, while the ϵ parameter in PSFT is fixed at 0.28. More details can be found in Appendix A.1.1.

4.1.1 TRAINING DYNAMICS

We begin by examining the training dynamics of each method. Specifically, we train each for around 10 epochs and plot the corresponding changes in entropy and performance. In this section, we use AIME-24 avg@32 to measure in-domain performance and GPQA (Rein et al., 2024) avg@8 to assess out-of-domain performance.

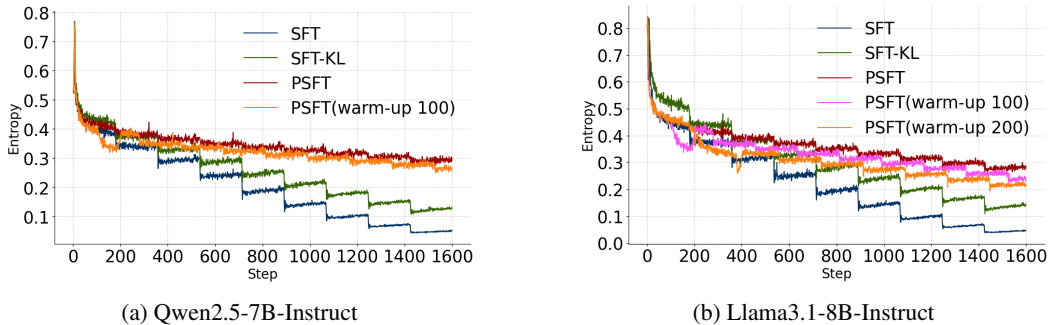


Figure 1: Training dynamics of Entropy. Each 178 steps is one epoch.

Observation 1: PSFT avoids entropy collapse. The entropy evolution is shown in Figure 1. Compared to SFT and SFT-KL, PSFT produces a smoother entropy curve. For SFT and SFT-KL, entropy exhibits a marked decline after each epoch, indicating potential overfitting. This suggests that PSFT is capable of sustaining long-term, token-level fine-grained training without triggering entropy collapse. Notably, PSFT with a warm-up phase demonstrates the same stability, and the number of warm-up rounds plays a critical role in shaping the overall entropy level. Further evidence is provided in Figure 6a.

Analysis 1: PSFT matches or surpasses in-domain performance of standard SFT at a similar entropy level. The in-domain evaluation performance evolution is shown in Figure 2. PSFT

²The dataset undergoes filtering through Math-Verify to exclude instances with unverifiable answers or responses exceeding 8,192 tokens in length.

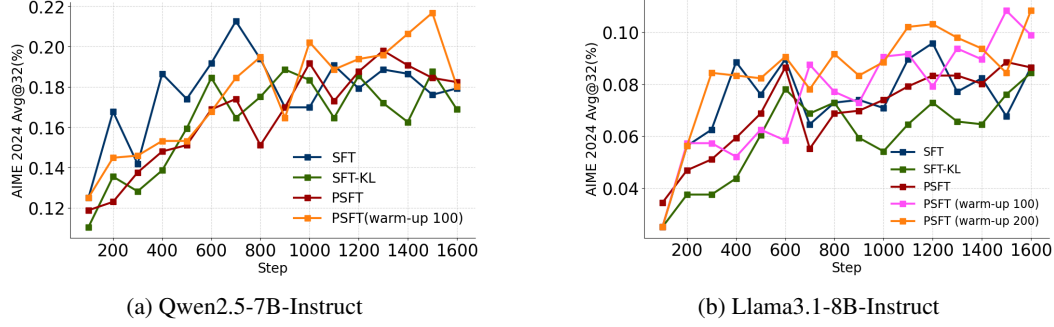


Figure 2: Training dynamics of in-domain performance.

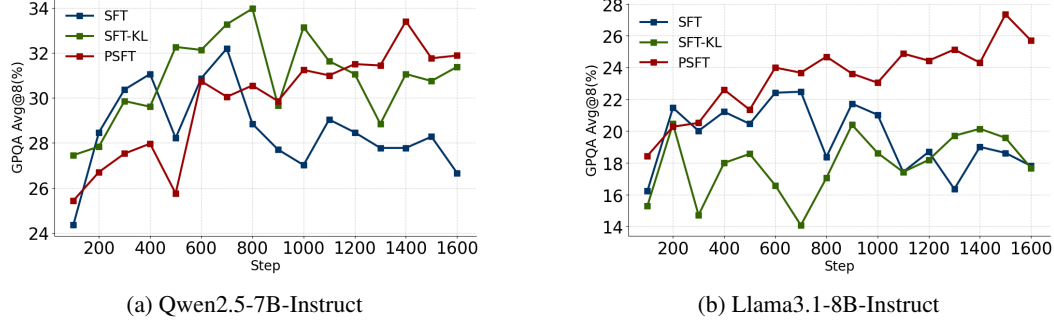


Figure 3: Training dynamics of out-of-domain performance.

achieves comparable performance on targeted math tasks and can even outperform SFT under the same entropy level. For instance, at step 1300, PSFT attains an AIME24 score of nearly 20 with an entropy around 0.3. This entropy level corresponds to that of SFT between epochs 300 and 520, yet PSFT achieves a higher AIME score than SFT within this range. While the KL-based approach intuitively addresses distribution shift, considering both entropy trends and in-domain performance, it remains less effective than PSFT.

Analysis 2: PSFT with warm-up operation surpasses standard SFT. Figure 2 illustrates the impact of introducing a warm-up phase into PSFT. We observe that the in-domain performance increases steadily. Across both models, PSFT with warm-up consistently outperforms vanilla PSFT as well as the standard SFT baselines. Moreover, increasing the number of warm-up steps further enhances in-domain performance. **This offers an effective way to fully exploit high-quality datasets** such as s1k (Muennighoff et al., 2025) and LIMO (Ye et al., 2025) while avoiding entropy collapse.

Analysis 3: PSFT can largely maintain the generalization of the model. The evolution of out-of-domain evaluation performance is shown in Figure 3. Increasing the number of training steps leads to a decline in generalization for SFT models. For example, after 700 steps, the SFT method on Qwen2.5-7B-Instruct exhibits a clear decline in performance. While SFT-KL alleviates this degradation to a certain degree, PSFT demonstrates substantially better generalization.

Overall, benefiting from its robust entropy control during training, **PSFT avoids large fluctuations in policy updates.** As observed in both in-domain and out-of-domain evaluations, PSFT generally maintains an upward performance trend. In contrast, SFT and SFT-KL exhibit pronounced fluctuations in out-of-domain performance.

4.1.2 DETAILED EVALUATIONS

Setup. (1) *Model setup:* For a fair comparison, we select checkpoints based on in-domain performance. For instance, on Qwen2.5-7B-Instruct, we use the SFT model trained for 700 steps, the SFT-KL model trained for 900 steps, and the PSFT model trained for 1300 steps. (2) *Evaluation*

Table 1: Detailed results of **in-domain** performance. For AIME and AMC, the results are avg.@32. For the rest, the results are avg.@8.

Method	AIME24	AIME25	AMC	MATH-500	OlympidBench	Minerva	Avg.
<i>Qwen2.5-7B-Instruct</i>							
Original	11.25	8.75	52.58	75.05	39.72	40.53	37.98
SFT	22.08	23.02	62.73	84.10	52.35	43.66	47.99
SFT _{KL}	19.27	21.56	63.20	83.55	52.70	42.19	47.08
PSFT _{warm-up}	22.92	23.02	62.42	84.68	52.30	43.38	48.17
PSFT	19.38	21.98	62.34	83.35	51.50	43.33	46.98
<i>Llama3.1-8B-Instruct</i>							
Original	3.96	0.73	24.45	48.55	17.09	25.87	20.11
SFT	10.63	16.77	47.19	72.60	41.43	32.17	36.80
SFT _{KL}	9.58	16.15	45.55	69.83	39.19	26.75	34.51
PSFT _{warm-up}	12.08	18.75	49.45	74.15	42.07	33.64	38.36
PSFT	10.31	14.48	46.80	71.98	39.61	32.40	35.93

Table 2: Detailed results of **out-of-domain** performance. For GPQA, ARC-C, TruthfulQA, and IFEval, the results are avg.@8. For the rest, the results are pass.@1.

Method	GPQA	ARC-C	TruthfulQA	MMLU-Pro	SuperGPQA	HeadQA	IFEval _{loose}	Avg.
<i>Qwen2.5-7B-Instruct</i>								
Original	31.38	91.54	66.10	54.99	27.59	73.41	73.94	59.85
SFT	32.89	92.22	63.14	58.98	29.02	74.65	54.42	57.90
SFT _{KL}	32.95	91.86	61.31	58.35	27.69	74.07	55.44	57.38
PSFT _{warm-up}	33.27	92.09	66.37	59.28	28.37	75.24	55.07	58.53
PSFT	33.21	92.29	67.16	59.18	28.10	75.82	73.03	61.26
<i>LLama3.1-8B-Instruct</i>								
Original	24.62	80.96	55.14	43.70	18.16	68.20	78.10	52.70
SFT	19.38	87.73	67.08	50.29	21.95	73.52	33.45	50.49
SFT _{KL}	18.18	87.02	67.03	47.02	19.81	71.95	34.19	49.31
PSFT _{warm-up}	23.99	89.72	68.55	56.03	25.26	77.71	33.08	53.48
PSFT	26.89	89.11	68.69	56.58	25.85	77.90	69.75	59.25

benchmark: (i) *In-domain tasks:* AIME24, AIME25, AMC, MATH-500 (Hendrycks et al., 2021), OlympidBench (He et al., 2024), and Minerva (Lewkowycz et al., 2022). (ii) *Out-of-domain tasks:* GPQA (Rein et al., 2024), ARC-C (Clark et al., 2018), TruthfulQA (Lin et al., 2021), MMLU-Pro (Wang et al., 2024), SuperGPQA (Du et al., 2025), HeadQA (Vilares & Gómez-Rodríguez, 2019) and IFEval (Zhou et al., 2023). (3) *Inference.* The inference length is set to 10,240 tokens (4,096 for IFEval), with a top- p of 0.95 and a temperature of 0.7.

Analysis 1. PSFT with warm-up consistently outperforms standard SFT on in-domain tasks.

Table 1 presents the performance across various in-domain benchmarks, where all fine-tuned methods achieve notable improvements over the original model. PSFT with a warm-up phase not only stabilizes early training dynamics as illustrated in Figure 2, but also enhances robustness across different evaluation benchmarks.

Analysis 2: PSFT demonstrates strong generalization. Table 2 shows different performance on out-of-domain benchmarks. Overall, our results are consistent with those reported in Zhou et al. (2025): injecting long CoT data into the model improves its general reasoning ability. However, on IFEval, which evaluates instruction-following capability, SFT, SFT-KL, and PSFT with warm-up significantly degrade the original performance. In contrast, PSFT not only preserves instruction compliance but also substantially enhances the model’s reasoning ability across other domains. PSFT with a warm-up phase also consistently surpasses standard SFT on out-of-domain tasks.

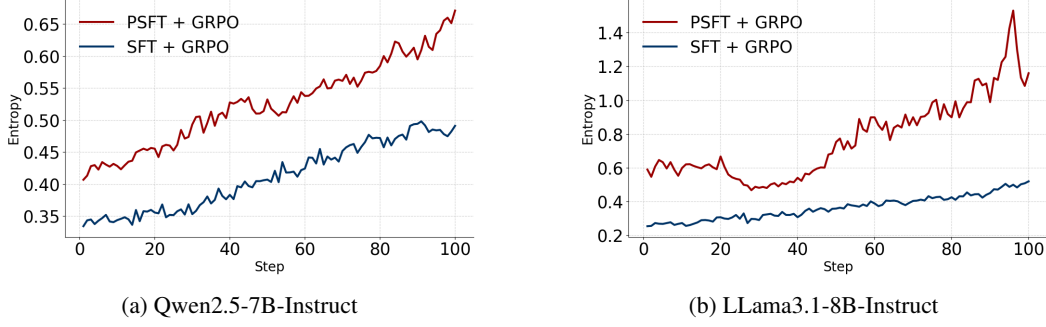


Figure 4: Training dynamics of Entropy on RL experiments.

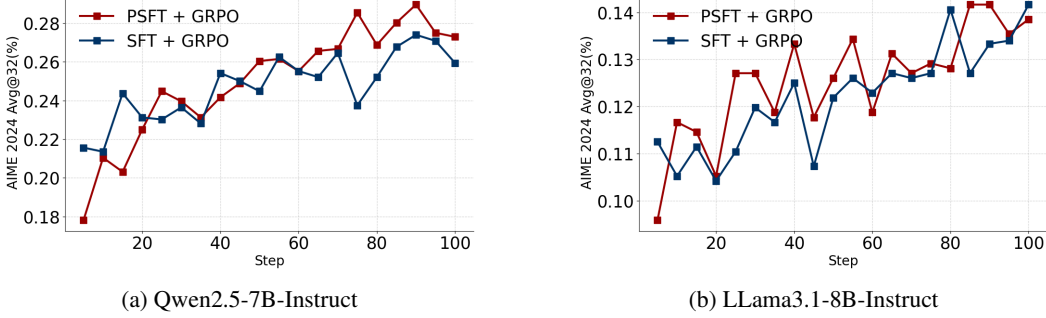


Figure 5: Training dynamics of in-domain performance on RL experiments.

Analysis 3: PSFT demonstrates robustness across models. For example, when evaluating GPQA and TruthfulQA tasks, SFT and SFT-KL exhibit varying behaviors across different models—one sometimes compromises performance while the other may boost it. In contrast, PSFT consistently improves performance, demonstrating a strong advantage over Llama models.

4.2 EXPLORATION ON THE POTENTIAL OF MODELS IN THE RL STAGE

Practical LLMs often conduct RL after SFT in post-training. The supervised fine-tuned models should function well as an appropriate start point of the following RL stage, avoid both overfitting and underfitting, and better stimulate the potential of RL. In this subsection, we evaluate the power of PSFT in the RL stage.

Setup. We adopt all techniques from DAPO (Yu et al., 2025), setting clip-higher to 0.28. The RL training uses the DAPO-MATH-17k dataset (Shao et al., 2024), with detailed training configurations provided in Appendix A.1.2. We further select the checkpoints based on the highest in-domain performance for evaluation.

4.2.1 TRAINING DYNAMICS ON RL

Analysis 1: PSFT leaves much room for RL optimization. The RL entropy evolution is shown in Figure 4. As low entropy limits the model’s exploration, it is notable that using PSFT as a cold start results in higher entropy throughout training, with a steeper increase.

Analysis 2: PSFT shows slow start and rapid catch-up in RL. The in-domain performance evolution is presented in Figure 5. Thanks to its high entropy, which promotes exploration, using PSFT as the cold start leads to a relatively steep performance improvement and ultimately surpasses the performance achieved when using SFT as the cold start.

Table 3: Detailed results of **in-domain** performance on RL experiments. For AIME and AMC, the results are avg.@32. For the rest, the results are avg.@8.

Method	AIME24	AIME25	AMC	MATH-500	OlympicBench	Minerva	Avg.
<i>Qwen2.5-7B-Instruct</i>							
SFT	22.08	23.02	62.73	84.10	52.35	43.66	47.99
↪ GRPO	27.40	26.15	70.86	86.60	58.09	45.27	52.40
PSFT	19.38	21.98	62.34	83.35	51.50	43.33	46.98
↪ GRPO	28.13	27.19	71.72	87.48	58.50	46.83	53.31
<i>Llama3.1-8B-Instruct</i>							
SFT	10.63	16.77	47.19	72.60	41.43	32.17	36.80
↪ GRPO	14.06	17.50	54.38	77.03	47.13	33.09	40.53
PSFT	10.31	14.48	46.80	71.98	39.61	32.40	35.93
↪ GRPO	14.27	18.02	55.23	77.78	47.74	37.55	41.77

Table 4: Detailed results of **out-of-domain** performance on RL experiments. For GPQA, ARC-C, TruthfulQA, and IFEval, the results are avg.@8. For the rest, the results are pass.@1.

Method	GPQA	ARC-C	TruthfulQA	MMLU-Pro	SuperGPQA	HeadQA	IFEval _{loose}	Avg.
<i>Qwen2.5-7B-Instruct</i>								
SFT	32.89	92.22	63.14	58.98	29.02	74.65	54.42	57.90
↪ GRPO	39.39	92.25	61.93	62.61	33.80	75.42	53.89	59.90
PSFT	33.21	92.29	67.16	59.18	28.10	75.82	73.03	61.26
↪ GRPO	43.43	92.42	64.84	63.65	34.51	75.89	73.73	64.06
<i>Llama3.1-8B-Instruct</i>								
SFT	19.38	87.73	67.08	50.29	21.95	73.52	33.45	50.49
↪ GRPO	31.06	88.57	64.82	54.39	27.80	73.81	31.85	53.19
PSFT	26.89	89.11	68.69	56.58	25.85	77.90	69.75	59.25
↪ GRPO	36.23	90.92	65.24	60.92	31.30	78.12	71.02	61.96

4.2.2 DETAILED EVALUATION

Analysis 3: PSFT achieves better in-domain performance after RL. Table 3 presents the results of RL experiments on in-domain benchmarks. Overall, RL improves the performance of all models. A comprehensive evaluation of the in-domain tasks shows that, although PSFT initially lags behind SFT, it holds significant potential that can be further unlocked through RL.

Analysis 4: PSFT outperforms SFT on out-of-domain RL tasks by a large margin. Table 4 presents the results of RL experiments on out-of-domain benchmarks. It is evident that if the original model exhibits weak capabilities during the cold-start phase, the subsequent RL training is constrained by these limitations. For instance, the IFEval results very clearly reflect this point. This further underscores the necessity of enhancing the current SFT algorithm.

4.3 FURTHER EXPERIMENTS ON HUMAN ALIGNMENT

In this section, to demonstrate the universality of PSFT, we extend our experiments to human alignment datasets, employing different base models and algorithms such as DPO (Rafailov et al., 2023). Our results show that PSFT effectively reduces the alignment tax (i.e., the generalization gap) while still leaving room for further optimization.

Setup. We adopt a completely different setup to demonstrate the universality of PSFT further. (1) *Model and Dataset.* We fine-tune the pre-trained Qwen3-4B-Base (Yang et al., 2025) model with the UltraFeedback dataset (Cui et al., 2023). (2) *Training algorithm.* We select the chosen region of the dataset and the SFT/PSFT model on it. And for PSFT, we train a double step named PSFT_{prolong}. Then, use the DPO algorithm to enable the model to learn the contrastive reward signal.

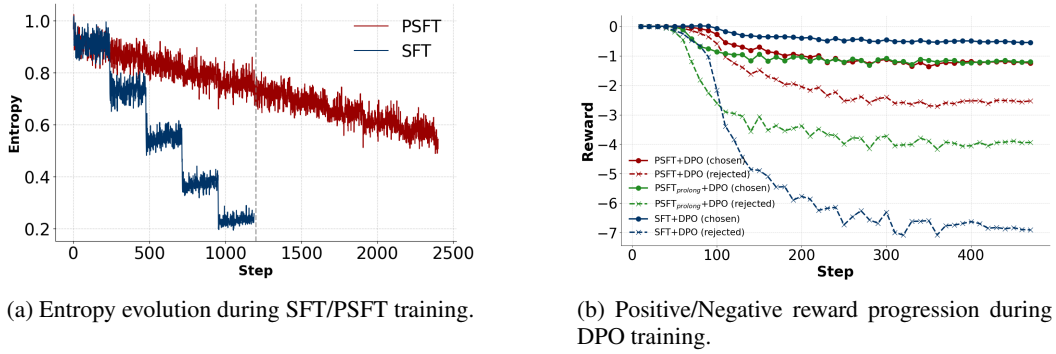


Figure 6: Training dynamics of SFT/PSFT followed by DPO

Table 5: Qwen3-4B-Base DPO training on the alignment benchmarks.

Method	AlpacaEval2		Arena-Hard WR (%)	MT-Bench	
	LC (%)	WR (%)		1-turn	2-turn
SFT	12.24	8.52	17.90	7.64	5.71
↔ DPO	16.96	13.40	26.50	7.91	6.00
PSFT _{prolong}	11.95	8.37	17.40	7.41	5.84
↔ DPO	19.26	15.17	30.20	7.63	6.74
PSFT	11.79	7.49	11.40	6.83	4.86
↔ DPO	23.29	20.13	36.40	8.51	6.95

(3) *Evaluation.* Since short CoT models can not solve the complex problems, we use the ARC, GSM8K (Cobbe et al., 2021), MMLU (Hendrycks et al., 2020), GPQA, and TruthfulQA to evaluate the alignment tax. We evaluate our models on three alignment benchmarks: MT-Bench (Zheng et al., 2024a), AlpacaEval (Dubois et al., 2024), and Arena-Hard (Li et al., 2024a). We use Qwen3-30B-A3-Instruct-2507 (Yang et al., 2025) as the judge model to provide alignment evaluation. We use llm-eval-harness (Gao et al., 2024) to evaluate the alignment tax performance. More experimental setup, see Appendix A.1.3.

4.3.1 TRAINING DYNAMICS

Analysis 1: PSFT still reliably prevents entropy collapse for other domains. As shown in Figure 6a, the entropy loss of SFT still shows a severe sawtooth shape, indicating a potential overfitting phenomenon, while PSFT overcomes this point. Therefore, SFT leads to a severe alignment tax (Figure 7) and restricts further optimization (Table 5).

4.3.2 HUMAN ALIGNMENT PERFORMANCE

Analysis2. Table 5 presents the results of DPO experiments on alignment benchmarks. We observe that PSFT_{prolong} achieves a similar effect to SFT training while avoiding entropy collapse, and it performs even better during DPO training. Moreover, using PSFT as the cold-start provides additional benefits that are further unlocked by DPO. This is consistent as reported by Xiao (2024) that when an SFT model collapses—i.e., becomes overly biased toward certain outputs with near certainty (Figure 6b)—it can subsequently induce preference collapse in the alignment process.

Analysis 3: PSFT training exploits both positive and negative samples well. As shown in Figure 6b, using SFT or PSFT as the cold-start point can influence the subsequent DPO phase. PSFT_{prolong} and PSFT yield similar rewards on positive samples; the main difference is that prolonged training reduces the occurrence of negative samples. In contrast, SFT training results in positive samples appearing frequently while negative samples are rare, which limits the model’s ability to learn from the reward signal and instead biases it toward the human-selected data.

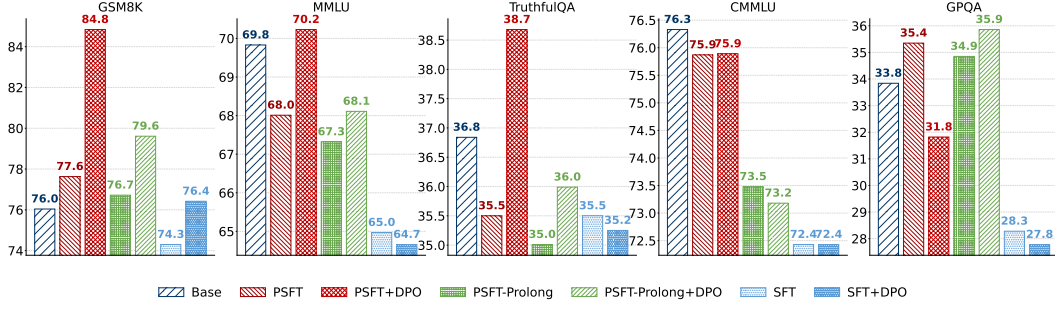
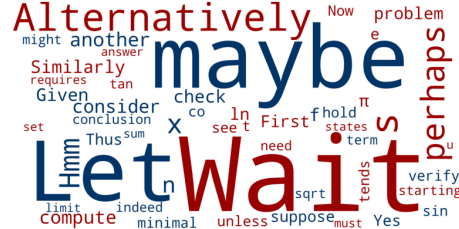


Figure 7: Results of models on alignment tax (out-of-domain tasks).



(a) Clipped token cloud in epoch 1.



(b) Clipped token cloud in epoch 3.

Figure 8: Examples and Changes of clipped tokens in PSFT during training.

Analysis4. Figure 7 shows the alignment tax results with various tasks (e.g., GSM8K, MMLU, GPQA) as the out-of-domain tasks. It indicates that: (a) PSFT and PSFT+DPO could largely maintain the general ability during the SFT stage compared to the original SFT. (b) PSFT-Prolong, even being degraded by possible overfitting issues, achieves relatively good results on alignment tax, implying the robustness of PSFT in practical scenarios.

5 IN-DEPTH ANALYSES

5.1 CLIPPED TOKENS IN PSFT

Figure 8 shows the representative clipped tokens during our PSFT’s training process. As we can see, the clipped token mostly concentrates on uncertain words like “wait”, “alternatively”, and similar words that reflect certain “long thinking patterns”, which are also supposed to be learned in the following RL stage. With the training in progress, these token-clip weights are more pronounced, while other tokens become smaller. The “thinking pattern” has been gradually and smoothly learned into models via PSFT, without much disturbance to the general capabilities.

5.2 PARAMETER ANALYSIS ON CLIPPED VALUE

As shown in Figure 9, PSFT without clipping maintains higher entropy but suffers from large, unstable gradient norms and highly fluctuating downstream results. By introducing the trust-region clipping mechanism, PSFT strikes a better balance: it prevents entropy collapse, stabilizes gradient updates, and achieves steady improvements in performance. Across different clipping thresholds, with medium values of ϵ (e.g., 0.28) showing particularly favorable trade-offs between stability and performance.

6 RELATED WORK

Supervised Fine-tuning. SFT is typically the initial step in the post-training pipeline. Addressing potential issues at the SFT stage can significantly enhance the robustness of the resulting further fine-tuned model. Several studies have noted that the standard cross-entropy (CE) loss may not be

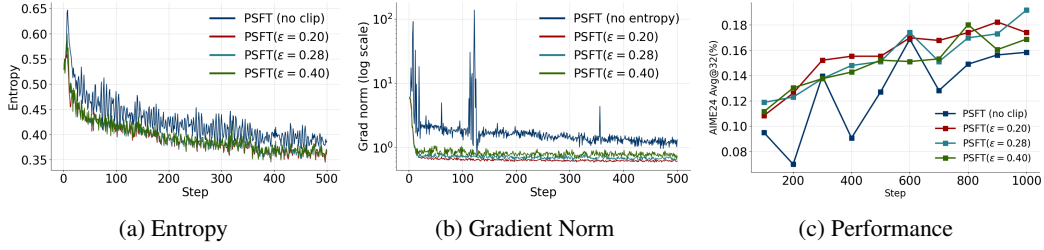


Figure 9: In-domain results of PSFT with different clipped values.

the most suitable objective for SFT (Li et al., 2024b; Xiao, 2024), arguing that training with CE at this stage tends to encourage the model to memorize the training data rather than acquire more generalizable capabilities. Li et al. (2024b) proposes the GEM, a game-theoretic SFT algorithm with entropy regularization to preserve diversity and mitigate forgetting. The concurrent work is Importance-weighted SFT (iw-SFT) (Qin & Springenberg, 2025) and DFT Wu et al. (2025). The key insight of iw-SFT is to reinterpret standard SFT as optimizing a loose lower bound on the RL objective in a sparse reward setting. While vanilla SFT simply maximizes the likelihood of filtered “successful” trajectories, this lower bound becomes increasingly loose as the model distribution drifts away from the reference policy. iw-SFT addresses this gap by introducing importance reweighting, which assigns higher weights to more preferred trajectories and thus tightens the bound toward the true RL objective. DFT views SFT as a flawed policy gradient and introduces a simple probability-based reweighting that improves generalization.

Reinforcement Learning. Following SFT, RL is employed to optimize reward signals directly. Policy gradient methods (Sutton et al., 2000) provide the foundation but are prone to instability and high variance. To address this, TRPO (Schulman et al., 2015) introduces a trust-region constraint, limiting each policy update to a small KL-divergence neighborhood to ensure stable improvement. PPO (Schulman et al., 2017) further simplifies this approach with a clipped surrogate objective, striking a balance between stability and efficiency.

7 CONCLUSION

Motivated by TRPO and PPO, we propose Proximal Supervised Fine-Tuning (PSFT). PSFT preserves a smooth entropy curve while achieving performance comparable to standard SFT at similar entropy levels, and it exhibits strong generalization. Furthermore, using PSFT as a cold-start model facilitates more effective post-training techniques such as RL optimization and DPO. Overall, PSFT presents a promising alternative to traditional SFT.

REFERENCES

- 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. *arXiv preprint arXiv:2501.17161*, 2025.
- 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.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback, 2023.

- Ganqu Cui, Yuchen Zhang, Jiacheng Chen, Lifan Yuan, Zhi Wang, Yuxin Zuo, Haozhan Li, Yuchen Fan, Huayu Chen, Weize Chen, et al. The entropy mechanism of reinforcement learning for reasoning language models. *arXiv preprint arXiv:2505.22617*, 2025.
- Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, King Zhu, Minghao Liu, Yiming Liang, Xiaolong Jin, Zhenlin Wei, et al. Supergpqa: Scaling llm evaluation across 285 graduate disciplines. *arXiv preprint arXiv:2502.14739*, 2025.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled alpaca-eval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.
- Hugging Face. Open r1: A fully open reproduction of deepseek-r1, 2025. URL <https://huggingface.co/blog/open-r1>.
- Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh, Nathan Lile, and Noah D Goodman. Cognitive behaviors that enable self-improving reasoners, or, four habits of highly effective stars. *arXiv preprint arXiv:2503.01307*, 2025.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model evaluation harness, 07 2024. URL <https://zenodo.org/records/12608602>.
- Etash Guha, Ryan Marten, Sedrick Keh, Negin Raoof, Georgios Smyrnis, Hritik Bansal, Marianna Nezhurina, Jean Mercat, Trung Vu, Zayne Sprague, et al. Openthoughts: Data recipes for reasoning models. *arXiv preprint arXiv:2506.04178*, 2025.
- 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.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng 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. *arXiv preprint arXiv:2402.14008*, 2024.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- 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.
- Maggie Huan, Yuetai Li, Tuney Zheng, Xiaoyu Xu, Seungone Kim, Minxin Du, Radha Pooven-dran, Graham Neubig, and Xiang Yue. Does math reasoning improve general llm capabilities? understanding transferability of llm reasoning. *arXiv preprint arXiv:2507.00432*, 2025.
- Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *Proceedings of the nineteenth international conference on machine learning*, pp. 267–274, 2002.
- 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.
- Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline. *arXiv preprint arXiv:2406.11939*, 2024a.

- Yang Li, Youssef Emad, Karthik Padthe, Jack Lanchantin, Weizhe Yuan, Thao Nguyen, Jason Weston, Shang-Wen Li, Dong Wang, Ilia Kulikov, et al. Naturalthoughts: Selecting and distilling reasoning traces for general reasoning tasks. *arXiv preprint arXiv:2507.01921*, 2025.
- Ziniu Li, Congliang Chen, Tian Xu, Zeyu Qin, Jiancong Xiao, Zhi-Quan Luo, and Ruoyu Sun. Preserving diversity in supervised fine-tuning of large language models. *arXiv preprint arXiv:2408.16673*, 2024b.
- Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- 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.
- OpenAI. Learning to reason with llms, 2024. URL <https://openai.com/index/learning-to-reason-with-llms/>.
- Chongli Qin and Jost Tobias Springenberg. Supervised fine tuning on curated data is reinforcement learning (and can be improved). *arXiv preprint arXiv:2507.12856*, 2025.
- 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.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Di-rani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.
- Yangjun Ruan, Neil Band, Chris J Maddison, and Tatsunori Hashimoto. Reasoning to learn from latent thoughts. *arXiv preprint arXiv:2503.18866*, 2025.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pp. 1889–1897. PMLR, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv: 2409.19256*, 2024.
- Richard S. Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 12, 2000.
- David Vilares and Carlos Gómez-Rodríguez. HEAD-QA: A healthcare dataset for complex reasoning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 960–966, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1092. URL <https://www.aclweb.org/anthology/P19-1092>.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil 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. *Advances in Neural Information Processing Systems*, 37:95266–95290, 2024.
- 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.

- Lechao Xiao. Rethinking conventional wisdom in machine learning: From generalization to scaling. *arXiv preprint arXiv:2409.15156*, 2024.
- Chulin Xie, Yangsibo Huang, Chiyuan Zhang, Da Yu, Xinyun Chen, Bill Yuchen Lin, Bo Li, Badih Ghazi, and Ravi Kumar. On memorization of large language models in logical reasoning. *arXiv preprint arXiv:2410.23123*, 2024.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- 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.
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024a.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand, 2024b. Association for Computational Linguistics. URL <http://arxiv.org/abs/2403.13372>.
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023. URL <https://arxiv.org/abs/2311.07911>.
- Ruochen Zhou, Minrui Xu, Shiqi Chen, Junteng Liu, Yunqi Li, Xinxin Lin, Zhengyu Chen, and Junxian He. Does learning mathematical problem-solving generalize to broader reasoning? *arXiv preprint arXiv:2507.04391*, 2025.

A APPENDIX

A.1 EXPERIMENTAL DETAILS

We perform SFT, PSFT, and RL training using the `verl` framework (Sheng et al., 2024), and employ `LLama-Factory` (Zheng et al., 2024b) for DPO training. The loss is aggregated using `token-mean` in `verl`. For SFT and PSFT, we use a weight decay of 0.1. All experiments are conducted with full fine-tuning.

A.1.1 MATH REASONING EXPERIMENTS

Method	Train batch size	Mini batch size	Learning rate	Train epochs	High clip	Cutoff_len
PSFT	256	32	1e-6	10	0.28	10k
SFT	256	–	2e-5	10	–	10k

Table 6: PSFT/SFT experiment configuration

The detailed training configurations for SFT/PSFT are presented in Table 6.

A.1.2 EXPLORATION ON MODEL POTENTIAL IN RL

The detailed training configurations for RL are presented in Table 7. For Qwen models and LLama models, we use the same configuration.

Config	SFT + GRPO	PSFT + GRPO
lr	1e-6	1e-6
kl_coef	0.0	0.0
max_prompt_length	2k	2k
max_response_length	10k	10k
overlong_buffer.len	2k	2k
train_batch_size	256	256
ppo_mini_batch_size	32	32
clip_ratio_low	0.2	0.2
clip_ratio_high	0.28	0.28
temperature	1.0	1.0
rollout.n	8	8
total_training_steps	100	100

Table 7: RL experiment configuration

Method	Train batch size	Mini batch size	Learning rate	Train epochs	High clip	Cutoff_len
PSFT	256	32	1e-6	5	0.28	6k
PSFT _{prolong}	256	32	1e-6	10	0.28	6k
SFT	256	–	2e-5	5	–	6k

Table 8: PSFT/SFT experiment configuration

Method	Train batch size	β	Learning rate	Train epochs	Cutoff_len
DPO	64	0.01	5e-7	1	4k

Table 9: DPO experiment configuration

A.1.3 ALIGNMENT EXPERIMENTS

Training. The detailed training configurations for SFT/PSFT on alignment tasks are presented in Table 8. The detailed training configurations for DPO are presented in Table 9.

Config	MT-Bench	Alpaca-Eval	Arena-Hard
temperature	0.0	0.7	1.0
top-p	–	0.95	0.7

Table 10: Alignment benchmark generation

Evaluation. We adopt the settings listed in Table 10 for evaluation generation. To assess alignment tax, we use the corresponding tasks in llm-eval-harness, as summarized in Table 11.

task	GSM8K	MMLU	TruthfulQA	CMMLU	GPQA
setting	gsm8k_cot	mmlu_flan_n_shot_generative	truthfulqa_mc1	cmmlu	cot_n_shot

Table 11: Alignment tax evaluation using llm-eval-harness

A.2 SYSTEM PROMPT

For the reasoning task, we use the system prompt as follows: Please reason step by step, and put your final answer within `\boxed{}`.

For the alignment task, we use the system prompt as follows: You are a helpful assistant.