

# BEYOND TWO-STAGE TRAINING: COOPERATIVE SFT AND RL FOR LLM REASONING

Liang Chen<sup>1</sup> Xueting Han<sup>2</sup> Li Shen<sup>3</sup> Jing Bai<sup>2</sup> Kam-Fai Wong<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong <sup>2</sup>Microsoft Research Asia

<sup>3</sup>Shenzhen Campus of Sun Yat-sen University

## ABSTRACT

Reinforcement learning (RL) has proven effective in incentivizing the reasoning abilities of large language models (LLMs), but suffers from severe efficiency challenges due to its trial-and-error nature. While the common practice employs supervised fine-tuning (SFT) as RL warmup, this decoupled two-stage approach suffers from catastrophic forgetting: second-stage RL gradually loses SFT-acquired behaviors and inefficiently explores new patterns. We introduce BRIDGE, a novel method to employ *bilevel optimization* to facilitate better cooperation between these training paradigms. By conditioning the SFT objective on the optimal RL policy, our approach enables SFT to meta-learn how to guide RL’s optimization process. During training, the lower-level performs RL updates while simultaneously receiving SFT supervision, while the upper-level explicitly maximizes the *cooperative gain*—the performance advantage of joint SFT-RL training over RL alone. Empirical evaluations across three LLMs and five reasoning benchmarks demonstrate that our method consistently outperforms baselines and achieves a better balance between effectiveness and efficiency. Specifically, BRIDGE achieves 44% faster training with a 13% performance gain on Qwen2.5-3B, and 14% faster training with a 10% improvement on Qwen3-8B.

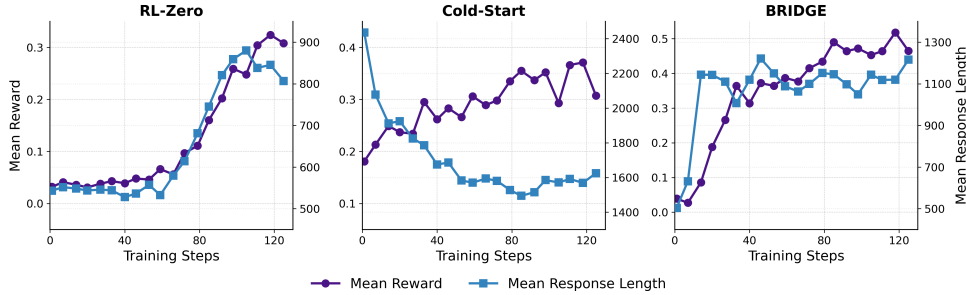


Figure 1: Training dynamics of mean reward and response length on Qwen2.5-3B.

## 1 INTRODUCTION

The emergence of OpenAI’s o1 (OpenAI) and DeepSeek-R1 (DeepSeek-AI et al., 2025) marks a significant advance in LLM reasoning capabilities, particularly for challenging tasks such as mathematics (Cobbe et al., 2021; Hendrycks et al., 2021b) and programming (Chen et al., 2021; Codeforces, 2025). The key technique driving this progress is large-scale, rule-based RL. However, the inherently trial-and-error nature of RL renders the training process highly inefficient (RL-zero in Figure 1). An alternative approach is SFT on curated long chain-of-thought (CoT) datasets, which enables models to rapidly acquire reasoning patterns through imitation learning. While more efficient, SFT typically exhibits poorer performance and generalization than RL (Chu et al., 2025).

\*Work was done during Liang Chen’s internship at MSRA. Contact: lchen@se.cuhk.edu.hk.

†Corresponding to: Xueting Han (chrihan@microsoft.com) and Kam-Fai Wong (kfwong@se.cuhk.edu.hk).

In practice, production-scale pipelines often use a two-stage cold-start paradigm, warming up with SFT before applying RL. This decoupled design introduces two issues: (a) catastrophic forgetting inherent to two-stage training, and (b) inefficient exploration in the second-stage RL due to the absence of guidance. As shown in Figure 1, the cold-start method exhibits a characteristic U-shaped trajectory in response length during the RL stage—first dipping and then rising—indicating that the model forgets previously learned expert patterns and explores new ones inefficiently. Meanwhile, the RL reward grows slowly due to the lack of guidance. As a result, the two stages fail to synergize, and both efficiency and performance remain unsatisfactory. This raises a natural question:

*Can we design a training framework that enables meaningful synergy between SFT and RL, ensuring their cooperation yields performance superior to standalone RL?*

To explore this possibility, we first propose a simple baseline that alternates between SFT and RL updates during training (Algorithm 1). Despite its simplicity, this approach improves both convergence efficiency and final performance (see Figure 2). However, such *independent updates* cannot guarantee improvements over RL alone, as not all SFT updates benefit RL optimization. Building on this insight, we develop BRIDGE, a cooperative learning framework based on bilevel optimization, where SFT serves as the upper-level problem and RL as the lower-level problem. By solving this nested structure—with the SFT objective explicitly conditioned on the RL solution—SFT provides *targeted* guidance that directly supports RL’s optimization process.

Specifically, BRIDGE employs an augmented model architecture comprising two learnable components: a base model and a LoRA module. The base model is optimized through the lower-level RL objective, while the LoRA parameters are updated via the upper-level supervised objective. To solve this bilevel problem, we adopt a first-order, penalty-based relaxation method. The relaxed lower-level update blends SFT and RL gradients, while the upper-level update *explicitly maximizes the cooperative gain—the performance advantage of joint SFT-RL training over RL-only optimization*. In this way, the lower level realizes the cooperation between two objectives, while the upper level ensures this cooperation yields superior performance.

To validate the effectiveness of our approach, we conduct experiments with three LLMs across five diverse benchmark datasets covering both standard and competition-level math reasoning tasks. Results demonstrate that BRIDGE consistently outperforms all baselines—including SFT, RL-zero, cold-start, and our alternating baseline—while requiring less wall-clock training time. These improvements confirm the benefits of tightly coupling SFT and RL through bilevel optimization rather than treating them as separate phases.

Our work makes the following contributions:

1. **Comparative analysis of reasoning training paradigms.** We systematically analyze three prevalent strategies for training large reasoning models. Our analysis reveals that the lack of interaction in two-stage pipelines prevents SFT and RL from effectively synergizing and leads to catastrophic forgetting and inefficient exploration. To mitigate these issues, we introduce a simple alternating baseline that achieves superior performance.
2. **A bilevel optimization framework for integrating SFT and RL.** To achieve deeper cooperation between SFT and RL, we propose BRIDGE, a bilevel optimization method that formalizes SFT as the upper-level and RL as the lower-level problem. Built on an augmented model architecture and solved using penalty-based relaxation, BRIDGE explicitly maximizes the cooperative gain—ensuring joint training outperforms standalone RL.
3. **Empirical validation on mathematical reasoning benchmarks.** We conduct extensive experiments with three LLMs across five mathematical reasoning benchmarks. BRIDGE consistently outperforms five baselines in both accuracy and training efficiency, demonstrating the practical benefits of tightly integrated SFT-RL optimization.

## 2 PRELIMINARIES

We begin by reviewing three prevalent fine-tuning strategies for training reasoning models, conduct a comparative analysis, and discuss limitations of the popular two-stage method. We then introduce a simple yet effective baseline that improves upon it.

## 2.1 FINE-TUNING METHODS FOR REASONING MODELS

We consider a language model parameterized by  $\theta$ , which defines a conditional distribution  $\pi(y|x; \theta)$  over output sequences  $y$  given input sequences  $x$ . This work focuses on three widely used methodologies for fine-tuning  $\theta$  to enhance the model’s reasoning capabilities.

**Supervised Fine-Tuning.** In supervised fine-tuning, we assume access to a curated dataset  $\mathcal{D}_{\text{SFT}} := \{(x, r, y)\}$  consisting of input prompts  $x$ , intermediate reasoning steps  $r$  distilled from larger reasoning models or annotated by human experts, and final answers  $y$ . The training objective maximizes the log-likelihood of generating both the reasoning process and the final answer:

$$\max_{\theta} J_{\text{SFT}}(\theta) := \mathbb{E}_{(x, r, y) \sim \mathcal{D}_{\text{SFT}}} [\log \pi(r, y | x; \theta)]. \quad (1)$$

This approach encourages the model to not only produce correct answers but also to imitate expert reasoning steps that lead to those answers.

**Reinforcement learning with verifiable rewards.** RLVR has gained increasing attention for its effectiveness in training advanced reasoning models such as DeepSeek-R1 (DeepSeek-AI et al., 2025). Given a dataset  $\mathcal{D}_{\text{RL}} := \{(x, y)\}$  with verifiable outputs—such as mathematics competition problems—the objective of rule-based RL is formulated as:

$$\begin{aligned} \max_{\theta} J_{\text{RL}}(\theta) := & \mathbb{E}_{(x, y) \sim \mathcal{D}_{\text{RL}}, (\hat{r}, \hat{y}) \sim \pi(\cdot | x; \theta)} [R(\hat{y}, y)] \\ & - \mathbb{E}_{(x, y) \sim \mathcal{D}_{\text{RL}}} [D_{\text{KL}}(\pi(\cdot | x; \theta) \parallel \pi_{\text{ref}}(\cdot | x))] \end{aligned} \quad (2)$$

where  $\pi_{\text{ref}}$  is a fixed reference model and  $R(\hat{y}, y)$  is a *rule-based reward function* that evaluates prediction correctness using a binary signal:

$$R(\hat{y}, y) = \begin{cases} 1, & \text{if } \hat{y} \equiv y, \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

Here,  $y$  denotes the ground-truth answer and  $\hat{y}$  is the model’s predicted output. The equivalence relation  $\hat{y} \equiv y$  is typically computed by a rule-based verifier. This objective is commonly solved using policy optimization methods such as Proximal Policy Optimization (PPO) (Schulman et al., 2017) or Group Relative Policy Optimization (GRPO) (DeepSeek-AI et al., 2025).

**Two-Stage Cold Start.** In practice, the common recipe uses SFT as a warm-up stage before applying RL. This two-stage approach, often referred to as "cold start," ensures that the model first learns to imitate expert reasoning patterns, providing a strong initialization for subsequent RL training.

## 2.2 COMPARISON OF FINE-TUNING METHODS

We evaluate these methods on mathematics problems at the grade 3–5 level. Figure 2 illustrates the evolution of test accuracy during training. We observe that *while SFT provides effective initialization and rapid early convergence for cold-start training, it contributes little to final convergence performance*. This results in faster initial accuracy improvements, but performance plateaus with minimal gains in the later stages of the two-phase pipeline. In contrast, RL alone converges more slowly but eventually achieves comparable final performance.

These results suggest that SFT and RL offer complementary strengths in reasoning tasks: SFT facilitates rapid initial learning, while RL enables better asymptotic performance. However, the naïve two-stage combination in cold-start training fails to fully exploit these complementary advantages. We identify two key limitations:

1. *Catastrophic forgetting:* The two-stage paradigm suffers from catastrophic forgetting—the model loses valuable SFT-acquired knowledge when transitioning to RL training. This phenomenon is evident in the response length dynamics during cold-start’s second stage (see the length dynamics in Figure 1). Response lengths initially drop sharply before gradually recovering, exhibiting a "dip-then-rise" pattern that indicates the model first forgets some expert behaviors before slowly exploring new strategies.
2. *Inefficient exploration:* Despite effective SFT initialization, online RL frequently encounters inefficient exploration, particularly on challenging problems where LLMs fail to generate reward-yielding solutions. LLMs often become trapped in local optima, unable to make further progress (see the reward dynamics in Figure 1). Moreover, once the initial SFT phase concludes, it cannot provide continued guidance for difficult problems.

These limitations motivate integrating SFT and RL training within a unified framework.

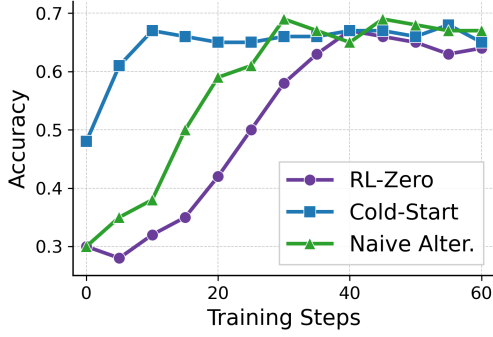


Figure 2: Comparison of Training Methods.

**Algorithm 1: A Simple Alternating Method**


---

```

1: Initialize parameters  $\theta_0$ ; datasets  $D_{\text{SFT}}$ ,  $D_{\text{RL}}$ ; learning rates  $\alpha_{\text{SFT}}$ ,  $\alpha_{\text{RL}}$ ; total steps  $T$ 
2: for  $t = 1$  to  $T$  do
3:   // RL step
4:   Sample  $(x_t, y_t) \sim D_{\text{RL}}$ 
5:   Generate solution with  $\pi_{\theta_{t-1}}(x_t)$ 
6:   Compute RL objective  $J_{\text{RL}}$  using (2)
7:    $\theta'_{t-1} \leftarrow \theta_{t-1} + \alpha_{\text{RL}} \nabla J_{\text{RL}}(\theta_{t-1})$ 
8:   // SFT step
9:   Sample example  $(x_t, r_t, y_t) \sim D_{\text{SFT}}$ 
10:  Compute SFT objective  $J_{\text{SFT}}$  using (1)
11:   $\theta_{t-1} \leftarrow \theta'_{t-1} + \alpha_{\text{SFT}} \nabla J_{\text{SFT}}(\theta'_{t-1})$ 
12: end for
    
```

---

### 2.3 A SIMPLE ALTERNATING BASELINE

To investigate the potential synergy between two methods, we design a simple alternating optimization strategy, as outlined in Algorithm 1. This approach alternates between RL steps, which explore novel reasoning strategies, and SFT steps, which imitate expert reasoning patterns.

As shown in Figure 2, this alternating strategy converges faster than pure RL and achieves better final performance than both standalone SFT and two-stage cold-start training. While this integration yields empirical gains, the current formulation treats SFT and RL as *independent update* processes with *no guarantee* that alternating updates will consistently outperform RL method alone. This limitation prompts us to investigate: *Can we develop a principled optimization framework where SFT and RL updates are inherently coordinated to guarantee improvements over pure RL?*

## 3 METHODOLOGY

In this section, we propose BRIDGE, a framework that tightly couples SFT and RL through a cooperative meta-learning approach. We first introduce the mathematical formulation, then present the learning algorithm and explanations.

### 3.1 BRIDGE: COOPERATIVE META-LEARNING FOR SFT AND RL

Given an SFT dataset  $\mathcal{D}_{\text{SFT}}$  and an RL dataset  $\mathcal{D}_{\text{RL}}$  (defined in Section 2.1), our objective is to integrate policy optimization (Eq. equation 2) with supervised learning (Eq. equation 1). We propose the following cooperative meta-learning formulation:

$$\begin{aligned}
 \max_w \quad & J_{\text{SFT}}(\theta^*(w), w) := \mathbb{E}_{(x,r,y) \sim \mathcal{D}_{\text{SFT}}} [\log \pi(r, y | x; \theta^*(w), w)] \\
 \text{s.t.} \quad & \theta^*(w) := \arg \max_{\theta} \left\{ \mathbb{E}_{(x,y) \sim \mathcal{D}_{\text{RL}}, (\hat{r}, \hat{y}) \sim \pi(\cdot | x; \theta, w)} [R(\hat{y}, y)] \right. \\
 & \quad \left. - \mathbb{E}_{(x,y) \sim \mathcal{D}_{\text{RL}}} [D_{\text{KL}}(\pi(\cdot | x; \theta, w) \| \pi_{\text{ref}}(\cdot | x))] \right\}.
 \end{aligned} \tag{4}$$

where  $\theta$  denotes the base model parameters and  $w$  represents the Low-Rank Adaptation (LoRA) weights (Hu et al., 2021). Together, they form an augmented model with parameters  $\bar{\theta} := [\theta, w]$ .

For clarity, we express Equation equation 4 in simplified notation:

$$\begin{aligned}
 \max_w \quad & J_{\text{SFT}}(w, \theta^*(w)), \\
 \text{s.t.} \quad & \theta^*(w) := \arg \max_{\theta} J_{\text{RL}}(\theta, w).
 \end{aligned} \tag{5}$$

This formulation exhibits a bilevel optimization structure inspired by the leader-follower game. SFT acts as the leader (teacher) with access to the RL follower’s (student’s) optimal response  $\theta^*(w)$ , enabling it to provide targeted guidance. Conversely, RL optimizes the base parameters  $\theta$  given the auxiliary support from SFT through  $w$ . During training, these components interact dynamically, resulting in better cooperation. As illustrated in Figure 3, this structure enables

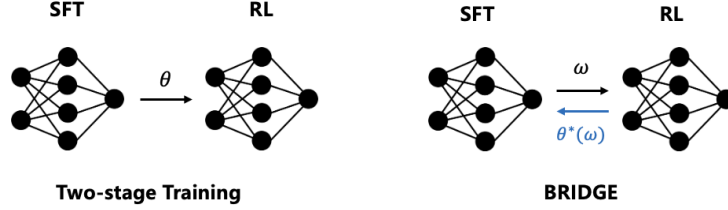


Figure 3: Comparison of two training methods.

*bidirectional* information flow—where RL’s optimal solution becomes visible to SFT—in contrast to the *unidirectional* flow of traditional two-stage approaches.

From a *meta-learning* perspective, BRIDGE implements cooperative framework where, at each iteration, the upper-level SFT provides an improved initialization for RL exploration, while the lower-level RL refines this initialization through reward-based optimization. This framework *adaptively extracts the most beneficial information from SFT to enhance RL training*, as SFT guidance may not always be uniformly beneficial.

The single-stage cooperative meta-learning design provides three advantages: (1) it avoids catastrophic forgetting of the two-stage pipeline through unified single-stage training; (2) it improves exploration efficiency via continuous supervised guidance; and (3) it guarantees RL performance gains by enabling SFT to meta-learn how to guide RL, strategically transferring beneficial knowledge.

**Architectural Design Rationale.** The augmented model architecture, comprising base model parameters  $\theta$  and LoRA parameters  $w$ , is essential for enabling cooperative learning. This separation allows the upper- and lower-level objectives to *co-adapt* during training, as illustrated in Figure 3. Without this architectural separation, our formulation (Equation equation 4) would collapse to a Model-Agnostic Meta-Learning (MAML)-style setup (Finn et al., 2017), where the lower-level solution reduces to a single gradient step used to update the upper-level SFT parameters. In this case, RL learning is disabled, and the cooperation between SFT and RL is lost.

### 3.2 LEARNING ALGORITHM

To solve the bilevel optimization problem in Eq. equation 5, we employ penalty-based methods (Shen & Chen, 2023; Shen et al., 2025) to avoid expensive second-order derivative computations. We first reformulate equation 5 as a single-level problem amenable to efficient first-order optimization.

We define the penalty function measuring the sub-optimality of the lower-level problem as:

$$p(w, \theta) = \max_{\theta'} J_{\text{RL}}(\theta', w) - J_{\text{RL}}(\theta, w). \quad (6)$$

This penalty quantifies the optimality gap:  $p(w, \theta) = 0$  if and only if  $\theta$  maximizes  $J_{\text{RL}}(\cdot, w)$ .

Given a penalty weight  $\lambda \in (0, 1)$ , we obtain the penalized reformulation:

$$\max_{\theta, w} \mathcal{L}(\theta, w) := (1 - \lambda) J_{\text{SFT}}(\theta, w) - \lambda p(w, \theta). \quad (7)$$

The penalty weight  $\lambda$  follows an annealing schedule: starting from a small value to warm-start training on supervised data, then gradually increasing to enforce the bilevel constraint more strictly.

Since  $\max_{\theta'} J_{\text{RL}}(\theta', w)$  depends only on  $w$ , the gradient with respect to  $\theta$  simplifies to:

$$\theta^{k+1} = \theta^k + \alpha [(1 - \lambda) \nabla_{\theta} J_{\text{SFT}}(\theta, w) + \lambda \nabla_{\theta} J_{\text{RL}}(\theta, w)] \quad (8)$$

For the gradient with respect to  $w$ , we invoke Danskin’s theorem. Assuming  $J_{\text{RL}}(\cdot, w)$  satisfies the required regularity conditions, we have:

$$\nabla_w \max_{\theta'} J_{\text{RL}}(\theta', w) = \nabla_w J_{\text{RL}}(\theta^*(w), w), \quad (9)$$

where  $\theta^*(w) = \arg \max_{\theta} J_{\text{RL}}(\theta, w)$ . In practice, we approximate  $\theta^*(w)$  by taking a single gradient ascent step with respect to the RL objective:

$$\hat{\theta} = \theta + \alpha \nabla_{\theta} J_{\text{RL}}(\theta, w), \quad (10)$$

**Algorithm 2:** Learning Algorithm of BRIDGE

---

```

1: Initialize augmented parameters  $\bar{\theta}^0 = (\theta^0, w^0)$ , and auxiliary parameters  $\hat{\theta}^0 := \theta^0$ ;
   learning rates  $\alpha, \beta$ ; penalty weight  $\lambda$ ; number of iterations  $K$ 
2: for  $k = 0$  to  $K - 1$  do
3:   Sample mini-batches  $\mathcal{B}_{\text{SFT}} \sim \mathcal{D}_{\text{SFT}}$  and  $\mathcal{B}_{\text{RL}} \sim \mathcal{D}_{\text{RL}}$ 
4:   // Compute base objectives
5:   Compute  $J_{\text{SFT}}(\theta^k, w^k)$ ,  $J_{\text{RL}}(\theta^k, w^k)$  and  $J_{\text{RL}}(\hat{\theta}^k, w^k)$  on  $\mathcal{B}_{\text{SFT}}$  and  $\mathcal{B}_{\text{RL}}$ 
6:   // Define composite objectives
7:    $J_{\text{Joint}}(\theta^k, w^k) = (1 - \lambda)J_{\text{SFT}}(\theta^k, w^k) + \lambda J_{\text{RL}}(\theta^k, w^k)$ 
8:    $J_{\text{Gain}}(w^k) = (1 - \lambda)J_{\text{SFT}}(\theta^k, w^k) + \lambda[J_{\text{RL}}(\theta^k, w^k) - J_{\text{RL}}(\hat{\theta}^k, w^k)]$ 
9:   // Update base parameters via joint objective
10:   $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta} J_{\text{Joint}}(\theta^k, w^k)$ 
11:  // Update auxiliary parameters via pure RL
12:   $\hat{\theta}^{k+1} \leftarrow \hat{\theta}^k + \alpha \nabla_{\hat{\theta}} J_{\text{RL}}(\hat{\theta}^k, w^k)$ 
13:  // Update LoRA parameters to maximize cooperative gain
14:   $w^{k+1} \leftarrow w^k + \beta \nabla_w J_{\text{Gain}}(w^k)$ 
15: end for

```

---

yielding the approximate gradient update for  $w$ :

$$\nabla_w \mathcal{L}(\theta, w) \approx (1 - \lambda) \nabla_w J_{\text{SFT}}(\theta, w) + \lambda \left[ \nabla_w J_{\text{RL}}(\theta, w) - \nabla_w J_{\text{RL}}(\hat{\theta}, w) \right]. \quad (11)$$

Algorithm 2 presents the learning procedure for BRIDGE. At each iteration, we sample SFT and RL mini-batches; update the base parameters  $\theta$  with the joint objective; update the auxiliary parameters  $\hat{\theta}$  via pure RL to track a baseline; and optimize the LoRA parameters  $w$  to maximize cooperative gain—the improvement of the joint objective over pure RL. This process lets SFT meta-learn how to guide RL’s optimization, mitigating catastrophic forgetting and inefficient exploration.

### 3.3 INTUITION BEHIND THE UPDATE RULES

**Lower-level update: Gradient fusion.** The update rule for  $\theta$  in Eq. equation 8 performs a convex combination of SFT and RL gradients. As  $\lambda$  increases from 0 to 1 during training, the algorithm smoothly transitions from imitation learning to reinforcement learning.

**Upper-level update: Maximizing cooperative gain.** The update for  $w$  in Eq. equation 11 solves the bilevel problem by finding LoRA parameters  $w$  that ensure the RL-optimized model  $\theta^*(w)$  also excels on the supervised dataset  $\mathcal{D}_{\text{SFT}}$ .

The update in Eq. equation 11 can be interpreted as performing gradient ascent on the following objective:

$$f(\theta, w) = (1 - \lambda) \underbrace{J_{\text{SFT}}(\theta, w)}_{\uparrow \text{likelihood on expert data}} + \lambda \underbrace{\left[ J_{\text{RL}}(\theta, w) - J_{\text{RL}}(\hat{\theta}, w) \right]}_{\uparrow \text{cooperative gain: SFT-RL vs RL-only}} \quad (12)$$

The first term maintains alignment with expert reasoning patterns, while the second term—the *cooperative advantage*—quantifies how much the joint SFT-RL optimization (using  $\theta$ ) outperforms pure RL training (using  $\hat{\theta}$ ). By maximizing this advantage term, the algorithm *explicitly encourages cooperation between supervised and reinforcement learning, ensuring their combination yields superior performance compared to RL alone*.

## 4 EXPERIMENT

### 4.1 SETTINGS

**Datasets.** We use LIMR (Li et al., 2025) (1.3k problems; Qwen2.5-3B) and MATH (Hendrycks et al., 2021a) (8.5k; lama-3.2-3B-Instruct and Qwen3-8B) for RL training. For the SFT dataset, we pair queries from LIMR and MATH with corresponding intermediate reasoning traces extracted

Table 1: Performance of BRIDGE compared to baselines on Qwen2.5-3B across five math benchmarks. Average relative performance gains (%) over Cold-start are highlighted in blue.

Method	MATH 500	Minerva Math	Olympiad Bench	AIME24	AMC23	Average
Base	32.4	11.8	7.9	0.0	20.0	14.4
SFT	53.4	18.8	21.5	3.3	42.5	27.9
RL-zero	64.4	26.5	27.0	3.3	40.0	32.2
Cold-start	66.0	24.3	26.8	9.0	35.0	32.2
Naive Alter.	65.2	25.3	27.1	6.7	42.5	33.4 (+3.7%)
BRIDGE	66.2	23.9	28.9	13.3	47.5	36.0 (+11.8%)

from DeepMath-103k (He et al., 2025), which were distilled from the DeepSeek-R1 model. We evaluate on five mathematical reasoning benchmarks: MATH500 (Hendrycks et al., 2021a), Minerva Math (Lewkowycz et al., 2022), OlympiadBench (He et al., 2024), AIME 2024, and AMC 2023, and two out-of-distribution (OOD) benchmarks (LiveCodeBench v5 (Jain et al., 2024) and GPQA (Rein et al., 2024)) to assess generalization.

**Models.** To demonstrate the generality of our approach, we experiment with three LLMs: Qwen2.5-3B (Yang et al., 2024), Llama-3.2-3B-Instruct (Grattafiori et al., 2024), and Qwen3-8B-Base (Yang et al., 2025). All models use prompt formats consistent with SimpleRL (Zeng et al., 2025). We use Verl (Sheng et al., 2024) for RL training; full settings are in Appendix B.

#### 4.2 BASELINES

We compare BRIDGE against five baselines on the same base models:

**Base/Instruction Model.** Base or instruction-tuned model with no reasoning-specific training.

**Supervised Fine-Tuning (SFT).** Trained only on curated reasoning traces by SFT (no RL).

**RL-Zero.** RL applied to the base model from scratch (no SFT warm-up).

**Cold-Start** Two-stage SFT then RL, with fully decoupled objectives.

**Naive Alternating.** Alternates SFT and RL updates without a cooperative objective.

#### 4.3 EXPERIMENTAL RESULTS

Table 2: Performance on Llama3.2-3B-Instruct.

Method	MATH 500	Minerva Math	Olympiad Bench	AIME24	AMC23	Average
Instruct	38.0	14.3	13.0	13.3	25.0	20.7
SFT	38.4	10.3	11.9	27.5	3.3	18.3
RL-zero	48.6	15.1	17.8	10.0	17.5	21.8
Cold-start	45.0	11.8	12.0	3.3	22.5	18.9
Naive Alter.	49.8	17.6	17.2	20.0	0.0	20.9 (+10.6%)
BRIDGE	51.8	15.1	19.3	10.0	27.5	24.7 (+30.7%)

**Main results across three LLMs.** Across five math benchmarks and three LLMs (Tables 1, 2, 3), BRIDGE attains the highest average accuracy. Cold-start behaves inconsistently across backbones: on Llama-3.2-3B-Instruct it is lower than RL-zero on average, whereas on Qwen3-8B-Base it exceeds RL-zero; this suggests that the two stage SFT then RL pipeline can constrain subsequent exploration and is not reliably optimal. Naive Alternating exceeds RL-zero and Cold-start on Qwen2.5-3B and Qwen3-8B-Base, indicating that training SFT and RL at the same time helps, although the gains are limited. In contrast, BRIDGE consistently outperforms all methods on all three backbones, with average improvements over Cold-start of 11.8% on Qwen2.5-3B, 30.9% on Llama-3.2-3B-Instruct, and 9.7% on Qwen3-8B-Base. Unlike naive alternating, BRIDGE treats SFT and RL as a cooperative



Table 3: Performance on Qwen3-8B-Base.

Method	MATH 500	Minerva Math	Olympiad Bench	AIME24	AMC23	Average
Base	55.4	24.3	22.5	3.3	27.5	26.6
SFT	67.8	32.0	29.8	45.0	13.3	37.6
RL-zero	76.2	36.0	42.4	10.0	50.0	42.9
Cold-start	80.4	38.2	39.6	16.6	52.5	45.5
Naive Alter.	78.2	37.5	40.6	13.3	65.0	46.9 (+3.1%)
BRIDGE	79.0	39.7	44.0	16.7	70.0	49.9 (+9.7%)

objective and adapts the supervision so that it benefits RL objective; this coupling yields robust and transferable gains across instruction tuned and base models and remains effective at larger scales.

**Generalization to More Challenging Math Tasks** Baseline methods tend to yield larger improvements on relatively easier benchmarks but generalize poorly to more complex reasoning tasks. For example, the Cold-start underperforms RL-zero on Minerva Math, OlympiadBench, and AMC23, indicating that Cold-start can restrict exploration and transfer poorly to harder problems. While the Naive Alternative partially mitigates this issue, its gains remain limited. In contrast, BRIDGE achieves consistent and substantial improvements on the more challenging benchmarks. These results underscore BRIDGE’s superior generalizability in handling complex mathematical reasoning.

**Generalization to Out-of-Domain tasks.** To test whether BRIDGE’s benefits transfer beyond math into out-of-domain reasoning. We evaluate on two non-math OOD benchmarks, LiveCodeBench and GPQA based on Qwen3-8B, as seen in Table 4, BRIDGE attains the best performance, demonstrating that joint training yields consistent cross-domain gains.

Table 4: OOD evaluation on Qwen3-8B.

Method	LiveCodeBench	GPQA	Avg.
Base	32.95	32.32	32.64
Instruct	33.07	35.80	34.44
RL-zero	32.50	38.38	35.44
Cold-start	23.86	25.76	24.81
BRIDGE	34.55	42.93	38.74

Table 5: Average performance across epochs.

Method	Epoch=1	Epoch=3	Epoch=6	Avg.
SFT	24.1	26.5	27.9	26.2
RL-zero	14.8	17.5	32.2	21.5
Cold-start	33.4	28.5	32.2	31.4
Naive Alt.	13.0	30.8	33.4	25.7
BRIDGE	32.3	33.3	36.4	34.0

**Training Dynamics Analysis.** We analyze the dynamics of mean reward and response length during training for BRIDGE, Cold-start, and RL-Zero on Qwen2.5-3B. As shown in Figure 1, the three methods exhibit markedly different patterns. RL-Zero suffers from online RL’s sample inefficiency, showing slow growth in both response length and reward. Cold-start begins with extremely long responses due to SFT warm-up, *causing slow training* (see in Table 6), followed by a sharp decline and gradual recovery. This "dip-then-rise" pattern indicates the model initially loses expert behavior acquired during SFT, then slowly explores new strategies—a mismatch that contributes to training inefficiency. Despite starting with higher rewards, Cold-start’s second-phase RL lacks proper guidance, resulting in final rewards similar to RL-Zero. In contrast, BRIDGE benefits from continuous SFT guidance throughout training, enabling rapid reward growth that surpasses Cold-start and achieving superior convergence. These dynamics demonstrate that BRIDGE’s bilevel optimization enables more efficient policy learning through sustained and targeted expert guidance.

Performance on downstream tasks across training epochs shows the same pattern (Table 5): Cold-start learns quickly at the beginning but dips mid-training and finishes no better than RL-zero, indicating constrained exploration. BRIDGE starts strong and improves steadily to the best final performance. Naive Alternating narrows the gap mid-training yet remains below BRIDGE. RL-zero shows the poorest early-stage efficiency. (see Appendix C for details.)

**Cost-Benefit Analysis.** We evaluated the cost-performance trade-offs by measuring wall-clock training time, average GPU memory usage per device, and final convergence performance across two model scales: Qwen2.5-3B (4×A100-80GB) and Qwen3-8B-Base (8×MI300-192GB). As shown



in Table 6, Cold-start requires nearly 2x the training time of RL-zero, despite the short SFT stage. This overhead stems from long sequence lengths induced by the SFT stage (see Figure 1). BRIDGE achieved 44% and 14% time savings compared to Cold-start for the 3B and 8B models, respectively. Despite a modest 11% increase in memory usage for the larger model, BRIDGE consistently delivered superior performance improvements (13% for 3B and 9.7% for 8B models), demonstrating favorable cost-benefit trade-offs for practical deployment.

Table 6: Cost-performance analysis on Qwen2.5-3B and Qwen3-8B-Base

Metric	Qwen 2.5-3B			Qwen 3-8B-Base		
	RL-zero	Cold-start	BRIDGE	RL-zero	Cold-start	BRIDGE
Time (hr)	6.1	12.3	6.9	38.5	39.1	33.5
Mem. (GB)	52.2	45.9	59.3	50.7	60.8	67.4
Acc. (%)	32.2	32.2	36.4	42.9	45.5	49.9

**LoRA ablation.** We prove BRIDGE is insensitive to LoRA hyperparameters (Appendix D).

## 5 RELATED WORK

Recent progress has highlighted the critical role of reinforcement learning in enhancing the reasoning capabilities of large language models (OpenAI; DeepSeek-AI et al., 2025). Recent advances have been made in recipes for training reasoning models. SimpleRL (Zeng et al., 2025) observes that fine-tuning on short-CoT datasets can harm reasoning ability, while He et al. (2025) find that fine-tuning on long-CoT distilled data can improve the reasoning performance of smaller models—especially when used as a warm-up stage before RL training. In practice, two-stage pipelines that combine SFT and RL are commonly used to balance stability and performance. However, existing approaches often rely solely on supervised fine-tuning, which tends to generalize poorly, or on pure RL, which suffers from sample inefficiency and unstable optimization.

Recent efforts move beyond the decoupled “SFT then RL” recipe by mixing two objectives within one stage. Concurrent work CHORD (Zhang et al., 2025) integrates SFT into on-policy RL via a weighted-sum objective with a global weight and a token-level weight that up-weights uncertain expert tokens and down-weights large-divergence tokens. This stabilizes training compared with naïve loss addition, but the coupling is heuristic and offers no mechanism to ensure that the injected supervision is useful for the RL update. LUFFY (Yan et al., 2025) combines off-policy expert traces with on-policy rollouts using Mixed-Policy GRPO, but it requires the off-policy demonstrations to be paired with the same prompts as the on-policy data, limiting flexibility and data reuse. BRIDGE instead treats SFT–RL cooperation as a bilevel problem, meta-adapting the supervision to maximize the reward gain of joint training over RL alone and yields larger and more robust improvements than simple loss mixing. It offers a new perspective on integrating imitation and exploration for large reasoning models.

We provide an extended discussion of related work in Appendix E.

## 6 CONCLUSION

This work investigates how to effectively integrate supervised fine-tuning and reinforcement learning to improve the reasoning capabilities of LLMs. We begin by analyzing three widely used training paradigms and identify a key limitation of existing multi-stage pipelines: the lack of interaction between SFT and RL. To address this, we propose a simple alternating baseline and further introduce BRIDGE, a bilevel optimization framework that models SFT as the upper-level objective and RL as the lower-level objective. By employing a penalty-based relaxation, BRIDGE explicitly encourages joint training to outperform standalone RL, fostering tighter cooperation between the two learning paradigms. Empirical results on five mathematical reasoning benchmarks demonstrate that our method consistently outperforms strong baselines in both accuracy and training efficiency. These findings underscore the potential of bilevel optimization as a unifying framework for combining supervised and reward-driven learning in complex reasoning tasks.

## REFERENCES

- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Sang Keun Choe, Sanket Vaibhav Mehta, Hwijeen Ahn, Willie Neiswanger, Pengtao Xie, Emma Strubell, and Eric Xing. Making scalable meta learning practical. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=Xazhn0JoNx>.
- 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. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=dYur3yabMj>.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Codeforces. Codeforces - competitive programming platform, 2025. URL <https://codeforces.com/>. Accessed: 2025-03-18.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojuan Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1126–1135. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/finn17a.html>.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan,

Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collet, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpiere Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippas Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Sweet, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings,

Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabisa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.

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.

Zhiwei He, Tian Liang, Jiahao Xu, Qiuzhi Liu, Xingyu Chen, Yue Wang, Linfeng Song, Dian Yu, Zhenwen Liang, Wenxuan Wang, Zhuosheng Zhang, Rui Wang, Zhaopeng Tu, Haitao Mi, and Dong Yu. Deepmath-103k: A large-scale, challenging, decontaminated, and verifiable mathematical dataset for advancing reasoning, 2025. URL <https://arxiv.org/abs/2504.11456>.

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*, 2021a.

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. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021b.

M Hong, HT Wai, Z Wang, and Z Yang. A two-timescale framework for bilevel optimization: Complexity analysis and application to actor-critic, dec. 20. *arXiv preprint arXiv:2007.05170*, 2020.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.

Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free

- evaluation of large language models for code, 2024. URL <https://arxiv.org/abs/2403.07974>.
- Prashant Khanduri, Siliang Zeng, Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A near-optimal algorithm for stochastic bilevel optimization via double-momentum. In *Advances in neural information processing systems*, 2021.
- Jeongyeol Kwon, Dohyun Kwon, Steve Wright, and Robert Nowak. On penalty methods for nonconvex bilevel optimization and first-order stochastic approximation. *arXiv preprint arXiv:2309.01753*, 2023.
- 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.
- Jiaxiang Li, Siliang Zeng, Hoi-To Wai, Chenliang Li, Alfredo Garcia, and Mingyi Hong. Getting more juice out of the sft data: Reward learning from human demonstration improves sft for llm alignment, 2024. URL <https://arxiv.org/abs/2405.17888>.
- Xuefeng Li, Haoyang Zou, and Pengfei Liu. Limr: Less is more for rl scaling, 2025. URL <https://arxiv.org/abs/2502.11886>.
- Xinyu Lin, Wenjie Wang, Yongqi Li, Shuo Yang, Fuli Feng, Yinwei Wei, and Tat-Seng Chua. Data-efficient fine-tuning for llm-based recommendation. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR ’24*, pp. 365–374, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400704314. doi: 10.1145/3626772.3657807. URL <https://doi.org/10.1145/3626772.3657807>.
- Songtao Lu. Slm: A smoothed first-order lagrangian method for structured constrained nonconvex optimization. 2024.
- OpenAI. Learning to reason with llms. [urlhttps://openai.com/index/learning-to-reason-with-llms/](https://openai.com/index/learning-to-reason-with-llms/). Accessed: 15 March 2025.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=Ti67584b98>.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Han Shen and Tianyi Chen. A single-timescale analysis for stochastic approximation with multiple coupled sequences. 2022.
- Han Shen and Tianyi Chen. On penalty-based bilevel gradient descent method. In *International Conference on Machine Learning*, 2023.
- Han Shen, Zhuoran Yang, and Tianyi Chen. Principled penalty-based methods for bilevel reinforcement learning and rlhf. 2024.
- Han Shen, Pin-Yu Chen, Payel Das, and Tianyi Chen. SEAL: Safety-enhanced aligned LLM fine-tuning via bilevel data selection. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=VHguhvc0M5>.
- 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.
- Reza Shirkavand, Qi He, Peiran Yu, and Heng Huang. Bilevel zofo: Bridging parameter-efficient and zeroth-order techniques for efficient llm fine-tuning and meta-training, 2025. URL <https://arxiv.org/abs/2502.03604>.

Quan Xiao, Han Shen, Wotao Yin, and Tianyi Chen. Alternating implicit projected sgd and its efficient variants for equality-constrained bilevel optimization. 2023.

Jianhao Yan, Yafu Li, Zican Hu, Zhi Wang, Ganqu Cui, Xiaoye Qu, Yu Cheng, and Yue Zhang. Learning to reason under off-policy guidance, 2025. URL <https://arxiv.org/abs/2504.14945>.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.

Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the wild, 2025. URL <https://arxiv.org/abs/2503.18892>.

Wenhao Zhang, Yuexiang Xie, Yuchang Sun, Yanxi Chen, Guoyin Wang, Yaliang Li, Bolin Ding, and Jingren Zhou. On-policy rl meets off-policy experts: Harmonizing supervised fine-tuning and reinforcement learning via dynamic weighting, 2025. URL <https://arxiv.org/abs/2508.11408>.

Rank/ $\alpha$	MATH500	Minerva	OlympiadBench	AIME24	AMC23	Avg.
32/16	79.0	39.7	44.0	16.7	70.0	49.9
16/32	79.0	38.6	44.0	16.0	70.0	49.5

Table 7: LoRA sensitivity ablation on Qwen3-8B-Base.

## A THE USE OF LARGE LANGUAGE MODELS

In preparing this manuscript, we used a large language model (LLM) solely for polishing the writing style and improving the clarity of the manuscript. The LLM was not used for generating research ideas, designing experiments, conducting analyses, or deriving results. All scientific contributions, including the conceptualization, methodology, experiments, and conclusions, were developed entirely by the authors.

## B IMPLEMENTATION DETAILS

All models are trained using the Verl framework (Sheng et al., 2024). We use a prompt batch size of 64, mini-batch size of 64, and learning rate of  $5 \times 10^{-7}$ . For LoRA, we set both rank and  $\alpha$  to 16. The penalty weight  $\lambda$  is set to 0.5. We employ two configurations: (1) for 3B models: 5 rollouts per prompt with 3k maximum tokens; (2) for 8B models: 8 rollouts per prompt with 8k maximum tokens. During evaluation, we use greedy decoding (temperature 0) with a 5k or 8k token limit and report pass@1 accuracy. Experiments are conducted on 4×NVIDIA A100 GPUs (80GB) for 3B models and 8×AMD MI300 GPUs (192GB) for 8B models.

## C PERFORMANCE ON VARIED FINE-TUNING EPOCHS

We assess BRIDGE’s effectiveness across different fine-tuning epochs on Qwen2.5-3B using average performance across epochs as the metric. As shown in Table 5, BRIDGE achieves the highest average performance. Among the baselines, Cold-start yields the second-best trade-off. However, its performance becomes unstable as training progresses, eventually converging to the same final result as RL-zero. In contrast, BRIDGE demonstrates consistent improvement throughout training. Overall, nearly all hybrid baselines outperform RL-zero in terms of early-stage efficiency, highlighting the advantage of integrating supervised fine-tuning and reinforcement learning paradigms.

## D LORA HYPERPARAMETER SENSITIVITY

To ensure BRIDGE’s gains do not hinge on a particular LoRA setting and to assess robustness to LoRA. We ablate LoRA rank ( $R$ ) and  $\alpha$  on Qwen3-8B-Base, keeping all training settings fixed and only changing  $(R, \alpha)$ . Table 7 indicates nearly identical outcomes across configurations, confirming that BRIDGE is insensitive to the LoRA choice.

## E RELATED WORK

**Reinforcement Learning for Large Reasoning Models.** Recent progress has highlighted the critical role of reinforcement learning in enhancing the reasoning capabilities of large language models (OpenAI; DeepSeek-AI et al., 2025). DeepSeek-R1 introduced a simple yet effective rule-based reward model and demonstrated further gains through multiple rounds of supervised distillation and RL training. LIMR (Li et al., 2025) showed that complex reasoning behaviors can emerge from as few as one thousand curated examples from the MATH dataset (Hendrycks et al., 2021b).

In parallel, substantial advances have been made in training recipes for large reasoning models. Chu et al. (2025) compare SFT and RL for reasoning tasks and find that RL generalizes significantly better, whereas SFT is prone to overfitting. SimpleRL (Zeng et al., 2025) observes that fine-tuning on short-CoT datasets can harm reasoning ability, while He et al. (2025) find that fine-tuning on long-CoT distilled data can improve the reasoning performance of smaller models—especially when



used as a warm-up stage before RL training. In practice, two-stage pipelines that combine SFT and RL are commonly used to balance stability and performance. However, existing approaches often rely solely on supervised fine-tuning, which tends to generalize poorly, or on pure RL, which suffers from sample inefficiency and unstable optimization. In this work, we propose the first unified training framework that enables explicit interaction between SFT and RL via a bilevel optimization formulation. This approach offers a new perspective on integrating imitation and exploration for large reasoning models.

**Bilevel Optimization in LLMs.** Bilevel optimization (BLO) is a classical framework for modeling hierarchical learning problems, originating from Stackelberg leader-follower games. Two major classes of methods have been developed to solve BLO problems. Implicit gradient methods (Hong et al., 2020; Khanduri et al., 2021; Shen & Chen, 2022; Xiao et al., 2023) compute gradients through the lower-level problem using second-order derivatives. While theoretically robust, these methods are often computationally expensive and memory-prohibitive when applied to large-scale models such as LLMs. In contrast, penalty-based relaxation methods (Shen & Chen, 2023; Kwon et al., 2023; Shen et al., 2024; Lu, 2024) approximate the BLO formulation using only first-order gradients, making them substantially more scalable and thus better suited for LLM applications. Recent work has explored the use of bilevel optimization in LLMs for tasks such as data selection (Lin et al., 2024; Shen et al., 2025), inverse reinforcement learning (Li et al., 2024), and meta-learning (Choe et al., 2023; Shirkavand et al., 2025). To the best of our knowledge, our work is the first to cast reasoning-oriented LLM training as bilevel optimization, introducing a novel augmented model architecture for modeling and solving this problem. This provides a principled framework for integrating supervised and reinforcement learning, where SFT actively assists RL optimization rather than merely serving as warmup.