

# NFT: BRIDGING SUPERVISED LEARNING AND REINFORCEMENT LEARNING IN MATH REASONING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Reinforcement Learning (RL) has played a central role in the recent surge of LLMs’ math abilities by enabling verification-driven training through binary verifier signals. In contrast, Supervised Learning (SL) is rarely considered for such verification-driven training, largely due to its heavy reliance on reference answers and inability to reflect on mistakes. In this work, we challenge the prevailing notion that self-improvement is exclusive to RL and propose Negative-aware Fine-Tuning (NFT) — a supervised approach that enables LLMs to reflect on their failures and improve autonomously with no external teachers. In online training, instead of throwing away self-generated negative answers, NFT constructs an *implicit* negative policy to model them. This implicit policy is parameterized with the same positive LLM we target to optimize on positive data, enabling direct policy optimization on all LLMs’ generations. We conduct experiments on 7B and 32B models in math reasoning tasks. Results consistently show that through the additional leverage of negative feedback, NFT significantly improves over SL baselines like rejection fine-tuning, matching, or even surpassing leading RL algorithms like GRPO and DAPO. Furthermore, we demonstrate that NFT and GRPO are actually equivalent in strict-on-policy training, even though they have entirely different theoretical foundations. Our experiments and theoretical findings bridge the gap between SL and RL methods in binary-feedback learning systems.

## 1 INTRODUCTION

The recent surge in math reasoning abilities of Large Language Models (LLMs) is largely driven by a fundamental shift in their learning paradigm, from imitation to self-improvement (DeepSeek-AI, 2025; OpenAI, 2024; Chen et al., 2025c). Instead of relying on reference answers supplied by human annotators or stronger models (Liu et al., 2024; OpenAI, 2023), the new paradigm requires only a question dataset with a binary verifier to judge the correctness of self-generated answers. By reflecting on their own generation, LLMs can improve autonomously. This approach not only eliminates the need for costly data annotation but also removes competence ceilings imposed by external teachers, offering a promising path toward general intelligence (DeepSeek-AI, 2025; NVIDIA, 2025).

Reinforcement Learning (RL) appears to be a natural fit for such verification-driven training. Specific algorithms like PPO (Schulman et al., 2017) and GRPO (Shao et al., 2024) are explicitly designed to maximize reward signals, which can conveniently take the form of a binary verifier outcome. In contrast, Supervised Learning (SL) is rarely considered for realizing such verification-driven training. A common view holds that SL is inherently designed to memorize the positive data, rendering it unsuitable for self-reflective learning from negative mistakes (Chu et al., 2025a).

*In this work, we challenge the prevailing notion that verification-driven training is exclusive to RL, and demonstrate it can be similarly achieved within the supervised learning paradigm.*

We start with a simple SL baseline: Rejection Fine-Tuning (RFT) (Yuan et al., 2023b; Dong et al., 2023). At each iteration, an LLM generates answers to questions. A verifier helps reject all negative answers. The remaining positive ones are compiled into a dataset to fine-tune the LLM itself in a supervised manner. RFT has been demonstrated effective (Bai et al., 2022; Yuan et al., 2023a; Song et al., 2023; Touvron et al., 2023; Xiong et al., 2025). However, it prevents any learning from negative feedback. LLMs are encouraged to reinforce what they already perform well, rather than reflect on their mistakes — an ability we believe critical for achieving general intelligence.

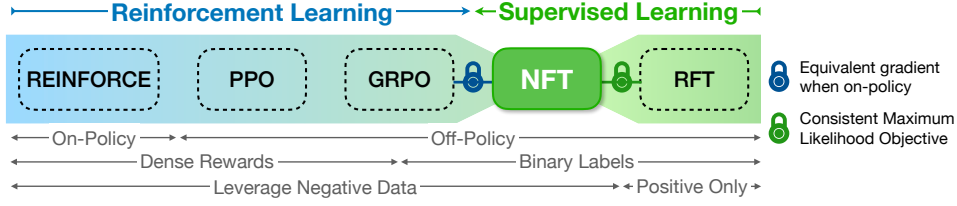


Figure 1: A spectrum of online algorithms for LLM fine-tuning. NFT bridges reinforcement learning and supervised learning methods through the leverage of negative feedback via supervision.

To overcome this limitation, we propose Negative-aware Fine-Tuning (NFT), an online learning algorithm that enables LLMs to learn from their negative generations (Sec. 3). Like RFT, NFT fine-tunes a positive LLM on positive answers via supervision. Crucially, instead of throwing away negative answers, NFT also constructs an *implicit* negative policy to model them. This implicit policy is parameterized with the same positive LLM we target to optimize on positive data, enabling direct policy optimization on all LLMs’ generations (Figure 2). NFT has minimal memory overhead, as only a single model is maintained throughout training.

To understand the connection between NFT and RL approaches, we conduct an in-depth comparison between NFT and GRPO (Sec. 4). Surprisingly, we find the two methods are actually equivalent in *strict-on-policy* training, despite that they originate from entirely different theoretical frameworks (Figure 1). Notably, the “advantage normalization” characteristic of GRPO is already implicitly reflected in NFT’s loss function. Their main difference arises in *off-policy* settings, regarding different strategies for clipping model gradients when the learned policy deviates from the old policy. These observations suggest a strong connection between SL and RL in binary-feedback learning systems.

We evaluate NFT on 7B and 32B Qwen models and report two key findings: 1. Supervised Learning alone can significantly enhance LLMs’ math reasoning with no external teachers. NFT matches or even surpasses state-of-the-art RL algorithms like GRPO (Shao et al., 2024) and DAPO (Yu et al., 2025). 2. The performance gap between RFT and RL in online training largely stems from SL’s past inability to leverage negative feedback, rather than any inherent superiority of RL. Through the additional leverage of negative data, NFT substantially mitigates this gap.

## 2 BACKGROUND

### 2.1 MAXIMUM LIKELIHOOD VS. POLICY GRADIENT

**Supervised Learning** essentially aims to learn a model  $\pi_\theta(\mathbf{a}|\mathbf{q})$  to fit the data distribution  $\pi(\mathbf{a}|\mathbf{q})$ . This can be realized by employing the maximum-likelihood objective:

$$\max_{\theta} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{q})} \log \pi_\theta(\mathbf{a}|\mathbf{q}) \Leftrightarrow \min_{\theta} D_{\text{KL}} [\pi(\mathbf{a}|\mathbf{q}) || \pi_\theta(\mathbf{a}|\mathbf{q})].$$

In LLM fine-tuning,  $\mathbf{q}$  usually means the question prompt, and  $\mathbf{a}$  the answer. To perform Maximum-likelihood training, we need a dataset  $\mathcal{D} = \{\mathbf{q}, \mathbf{a} \sim \pi(\mathbf{a}|\mathbf{q})\}$  to draw training samples from.

**Reinforcement Learning**, by contrast, maximizes pre-defined reward  $r(\mathbf{q}, \mathbf{a})$  for  $\mathbf{a} \sim \pi_\theta(\mathbf{a}|\mathbf{q})$ :

$$\max_{\theta} J(\theta) := \mathbb{E}_{\mathbf{a} \sim \pi_\theta(\mathbf{a}|\mathbf{q})} r(\mathbf{q}, \mathbf{a}).$$

Directly back-propagating through  $J(\theta)$  is non-trivial, as  $r(\cdot)$  can be an arbitrary scalar function whose gradient is unknown. Luckily,  $\nabla_{\theta} J(\theta)$  can be estimated, making policy optimization feasible.

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\mathbf{a} \sim \pi_\theta(\mathbf{a}|\mathbf{q})} \nabla_{\theta} [r(\mathbf{q}, \mathbf{a}) \log \pi_\theta(\mathbf{a}|\mathbf{q})]. \quad (1)$$

Eq. 1 is known as the Policy Gradient (PG) algorithm (Sutton et al., 1999). In sequential decision-making problems such as language reasoning,  $\mathbf{a}$  can be interpreted as the token decision for each step  $t$ , and  $r(\mathbf{q}, \mathbf{a})$  can be replaced with advantage functions  $A(\mathbf{q}, \mathbf{a})$  (Schulman et al., 2015).

### 2.2 MATH REASONING RL: FROM POLICY GRADIENT TO GRPO

Eq. 1 requires training to be *on-policy*, where  $\pi_\theta$  can only be updated a single time after data collection. To break this limitation, importance sampling can be applied (Schulman et al., 2017). Suppose

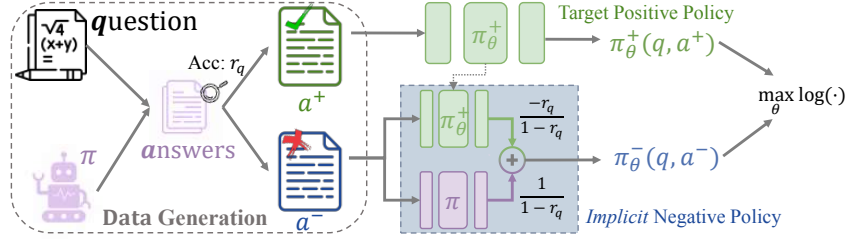


Figure 2: Illustration of the NFT algorithm. **Data Collection:** An LLM  $\pi$  generates answers to a set of math questions. Generation results are split into two sub-datasets based on their answer correctness. **Policy Optimization:** By constructing an *implicit* policy for modeling negative data, NFT enables direct policy optimization on both positive and negative answers via maximum-likelihood.

the policy for collecting the RL dataset is denoted as  $\pi_{\text{old}}$ , we have

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\mathbf{a} \sim \pi_{\text{old}}(\mathbf{a}|\mathbf{q})} \left[ \frac{\pi_{\theta}(\mathbf{a}|\mathbf{q})}{\pi_{\text{old}}(\mathbf{a}|\mathbf{q})} r(\mathbf{q}, \mathbf{a}) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}|\mathbf{q}) \right] = \mathbb{E}_{\mathbf{a} \sim \pi_{\text{old}}(\mathbf{a}|\mathbf{q})} [r(\mathbf{q}, \mathbf{a}) \nabla_{\theta} R_{\theta}(\mathbf{q}, \mathbf{a})], \quad (2)$$

where  $R_{\theta}(\mathbf{q}, \mathbf{a}) := \frac{\pi_{\theta}(\mathbf{a}|\mathbf{q})}{\pi_{\text{old}}(\mathbf{a}|\mathbf{q})}$  is the likelihood Ratio between two policies.

In math reasoning tasks, PPO (Schulman et al., 2017) and subsequent GRPO (Shao et al., 2024) algorithms further apply gradient clipping to constrain the distance between  $\pi_{\theta}$  and  $\pi_{\text{old}}$ :

$$\mathcal{L}^{\text{GRPO}}(\theta) = - \sum_{\mathbf{q}, \mathbf{a} \sim \pi_{\text{old}}} \sum_t \min \left[ R_{\theta}^t(\mathbf{q}, \mathbf{a}) \hat{A}_{\mathbf{q}, \mathbf{a}}, \text{clip}(R_{\theta}^t(\mathbf{q}, \mathbf{a}), 1 - \epsilon', 1 + \epsilon') \hat{A}_{\mathbf{q}, \mathbf{a}} \right], \quad (3)$$

where  $R_{\theta}^t(\mathbf{q}, \mathbf{a}) := \frac{\pi_{\theta}(\mathbf{a}_t|\mathbf{q}, \mathbf{a}_{<t})}{\pi_{\text{old}}(\mathbf{a}_t|\mathbf{q}, \mathbf{a}_{<t})}$ , and  $\hat{A}_{\mathbf{q}, \mathbf{a}}$  is the estimated advantage value. Note that we have dropped some auxiliary loss terms, such as KL and entropy regularization, as they have been pointed out to be unnecessary by recent studies like DAPO (Yu et al., 2025).

GRPO proposes an efficient way to estimate  $\hat{A}_{\mathbf{q}, \mathbf{a}}$ . Collect  $K$  answers  $\mathbf{a}^{1:K}$  and their binary reward  $r^{1:K} \in \{0, 1\}$  for each question. The advantage is defined to be the normalized reward:

$$\hat{A}_{\mathbf{q}, \mathbf{a}} := [r(\mathbf{q}, \mathbf{a}) - \text{mean}\{r^{1:K}\}] / \text{std}\{r^{1:K}\}. \quad (4)$$

Later studies (Liu et al., 2025b) suggest removing the  $\text{std}$  term from Eq. 4.

## 3 METHOD

### 3.1 PROBLEM SETUP

**Dataset.** Given a set of  $N$  math questions  $\{\mathbf{q}^{1:N}\}$ , a pretrained LLM  $\pi_{\text{old}}(\mathbf{a}|\mathbf{q})$ , and a verifier for judging the correctness. In every iteration, we generate a dataset  $\mathcal{D} = \{\mathbf{q}, \mathbf{a}^{1:K} \sim \pi, r^{1:K}\}^{1:N}$ , where  $r \in \{0, 1\}$  is the correctness label, and  $K$  the number of answers collected for each question.

Dataset  $\mathcal{D} \sim \pi$  can be split into two subsets:  $\mathcal{D}^+$  and  $\mathcal{D}^-$ .  $\mathcal{D}^+$  contains all positive answers, and  $\mathcal{D}^-$  contains the rest negative ones. We denote the underlying answer distribution of  $\mathcal{D}^+$  as  $\pi^+(\cdot|\mathbf{q})$ .

**Learning Target.** We want to optimize the old policy  $\pi$  into a new policy  $\pi_{\theta}^+ \approx \pi^+$ . The target  $\pi^+(\mathbf{a}|\mathbf{q})$  can be formalized using Bayes' Rule:

$$\pi^+(\mathbf{a}|\mathbf{q}) := \pi(\mathbf{a}|\mathbf{q}, r=1) = \frac{\pi_{\text{old}}(\mathbf{a}|\mathbf{q})p(r=1|\mathbf{q}, \mathbf{a})}{\sum_A \pi_{\text{old}}(\mathbf{a}|\mathbf{q})p(r=1|\mathbf{q}, \mathbf{a})}, \quad (5)$$

where  $A$  means all possible language space for  $\mathbf{a}$ .

**Discussion.** An obvious solution for learning  $\pi^+$  is to fine-tune solely on correct answers ( $\mathcal{D}^+$ ) and discard  $\mathcal{D}^-$  totally (RFT) (Dong et al., 2023; Xiong et al., 2025). However, this approach prevents the model from any learning on its negative feedback ( $\mathcal{D}^-$ ). We posit that the ability to reflect on one's own failures is not merely desirable, but central to general intelligence, marking a shift from

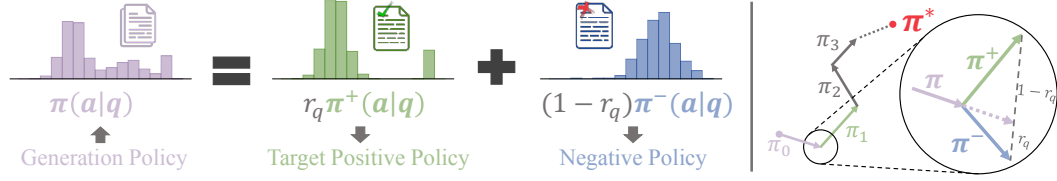


Figure 3: **Left:** Policy Splitting. The generation policy can be split into a positive policy and a negative policy, and re-expressed as their linear combination. **Right:** Policy Improvement. By iteratively optimizing towards its positive split, an LLM policy  $\pi_0$  can improve continuously.

pure imitation to self-reflective learning. Though traditionally viewed as a distinctive strength of RL (Chu et al., 2025a; Yue et al., 2025), we ask: *Can self-reflective improvement be similarly achieved within the SL paradigm?*

### 3.2 DIRECT OPTIMIZATION OF LANGUAGE MODELS WITH NEGATIVE ANSWERS

In this section, we discuss how to leverage negative data  $\mathcal{D}^-$  to directly optimize  $\pi_\theta^+$ . Despite seemingly impossible at first thought, we find the target policy  $\pi^+$  and the negative policy  $\pi^-$  are tightly coupled, making feasible training  $\pi_\theta^+$  directly from  $\mathcal{D}^-$ .

First, we formalize the definition of the negative policy  $\pi^-$  similar to Eq. 5

$$\pi^-(a|q) := \pi(a|q, r=0) = \frac{\pi_{\text{old}}(a|q)[1 - p(r=1|q, a)]}{\sum_A \pi_{\text{old}}(a|q)[1 - p(r=1|q, a)]}. \quad (6)$$

Combining Eq. 5 and Eq. 6, we make a key observation that

$$r_q \pi^+(a|q) + [1 - r_q] \pi^-(a|q) = \pi_{\text{old}}(a|q), \quad (7)$$

where  $r_q := \sum_A \pi_{\text{old}}(a|q)p(r=1|q, a) = p(r=1|q)$  is the correctness rate of LLM  $\pi_{\text{old}}$  over a question  $q$ . In practice,  $r_q \approx \text{mean}\{r^{1:K}\}$  can be estimated using the  $K$  rewards in dataset  $\mathcal{D}$ .

**Implicit negative policy.** Eq. 7 reveals a tight coupling between  $\pi^+$  and  $\pi^-$  (Figure 3). Given that we already have  $\pi$  as the pretrained LLM and  $r_q$  is estimable, learning  $\pi^-$  from negative data should, in principle, shape the target policy  $\pi_\theta^+$ , in a manner analogous to SL on positive data.

To realize this idea, we construct an *implicit* negative policy, denoted  $\pi_\theta^-$ , by re-parameterizing the target policy  $\pi_\theta^+$  using the relationship in Eq. 7:

$$\pi_\theta^-(a|q) := \frac{\pi_{\text{old}}(a|q) - r_q \pi_\theta^+(a|q)}{1 - r_q}.$$

Thus, training  $\pi_\theta^-$  on negative answers directly leads to optimizing the underlying positive LLM  $\pi_\theta^+$  (Figure 2). We have the following guarantee:

**Theorem 3.1 (Policy Optimization with Negative Answers).** *Consider the maximum-likelihood objective for training the implicit negative policy  $\pi_\theta^-$ :*

$$\max_{\theta} \mathbb{E}_{p(q)\pi^-(a|q)} [\log \pi_\theta^-(a|q)] \Leftrightarrow \min_{\theta} \left[ -\mathbb{E}_{(q,a) \sim \mathcal{D}^-} \log \frac{\pi_{\text{old}}(a|q) - r_q \pi_\theta^+(a|q)}{1 - r_q} \right] \quad (8)$$

Assuming unlimited data and model capacity, the optimal solution for solving Eq. 8 is

$$\forall q, a: \quad \pi_{\theta^*}^-(a|q) = \pi^-(a|q)$$

Proof in Appendix A. Theorem 3.1 demonstrates the feasibility of policy optimization on negative data only. To further utilize positive data, we additionally conduct supervised training on  $\mathcal{D}^+$ :

$$\mathcal{L}_{(a,q,r) \sim \mathcal{D}}^{\text{NFT}}(\theta) = r \left[ -\log \frac{\pi_\theta^+(a|q)}{\pi_{\text{old}}(a|q)} \right] + (1-r) \left[ -\log \frac{1 - r_q \frac{\pi_\theta^+(a|q)}{\pi_{\text{old}}(a|q)}}{1 - r_q} \right] \quad (9)$$

**Algorithm 1** Negative-aware Fine-Tuning (NFT)

---

```

1: Input: Language model  $\pi$ , prompt set  $\mathbf{q}^{1:N}$ , verifier  $r(\cdot)$ .
2: Def  $\max_v(\mathbf{x}, \epsilon)$ : //Straight-through Max Operator
3:   Return  $\text{stop\_gradient}[\max(\mathbf{x}, \epsilon) - \mathbf{x}] + \mathbf{x}$  //Clip Value while Keeping Gradient
4: for each iteration do
5:   for each sampled prompt  $\mathbf{q}$  do
6:     Generate  $K$  answers  $\mathbf{a}^{1:K}$  and verify their correctness  $r^{1:K}$  // Data Collection
7:     Calculate correctness rate  $\hat{r}_q = \text{mean}\{r^{1:K}\}$  and token-level likelihood  $\{\pi_{\text{old}}(\mathbf{a}_t|\mathbf{q}, \mathbf{a}_{<t})^{1:K}\}$ 
8:      $\mathcal{D} \leftarrow \{\mathbf{q}, \hat{r}_q, \mathbf{a}^{1:K}, r^{1:K}, \pi_t^{1:K}\}$  If  $0 < r_q < 1$  // Prompt Filtering
9:   end for
10:  Initialize  $\pi_\theta^+ \leftarrow \pi$ 
11:  for each mini batch  $\{\mathbf{q}, \mathbf{a}, r, \hat{r}_q, \pi_t\}$  in  $\mathcal{D}$  do
12:     $R_\theta^t(\mathbf{q}, \mathbf{a}) = \frac{\pi_\theta^+(\mathbf{a}_t|\mathbf{q}, \mathbf{a}_{<t})}{\pi_{\text{old}}(\mathbf{a}_t|\mathbf{q}, \mathbf{a}_{<t})}, \forall t$  // Positive Likelihood Ratio
13:    If  $r = 0$ :
14:       $R_\theta^t(\mathbf{q}, \mathbf{a}) = (1 - \hat{r}_q R_\theta^t(\mathbf{q}, \mathbf{a})) / (1 - \hat{r}_q)$  // Implicit Negative Likelihood Ratio
15:       $R_\theta^t(\mathbf{q}, \mathbf{a}) = \max_v[R_\theta^t(\mathbf{q}, \mathbf{a}), \epsilon]$  // Clip Negative Likelihood Ratio
16:       $\theta \leftarrow \theta + \lambda \nabla_\theta \sum_t \log R_\theta^t(\mathbf{q}, \mathbf{a})$  (Eq. 10) // Maximum Likelihood Training
17:    end for
18:     $\pi \leftarrow \pi_\theta^+$ , start the next iteration
19: end for

```

---

Note that we subtract a baseline term  $-\log \pi_{\text{old}}(\mathbf{a}|\mathbf{q})$  from the loss. Since this term is unrelated to  $\theta$ , it does not affect the loss gradient and thus the optimal solution.  $\pi_{\text{old}}(\mathbf{a}|\mathbf{q})$  is the old data likelihood before optimizer update. At the start of training, we have  $\pi_\theta^+ = \pi$  such that  $\mathcal{L}_{(\mathbf{a}, \mathbf{q}, r)}^\theta = 0$ .

We name our method Negative-aware Fine-Tuning (NFT) as it enables the additional leverage of negative data for policy optimization compared with RFT.

**Continuous Reward.** Though we mainly consider binary rewards for math reasoning, NFT is directly applicable to continuous rewards. Following Levine (2018), we can define optimality reward  $r \in [0, 1]$ , meaning the answer has a probability of  $r$  for being positive, and  $1 - r$  for being negative. The loss definition in Eq. 9 and its convergence property stay unchanged (Proof in Appendix B).

**Memory Efficiency.** NFT is memory-efficient. In practice, we keep only a single model copy in memory. The old policy likelihood  $\pi_{\text{old}}(\mathbf{a}|\mathbf{q})$  can be pre-computed during data generation.

### 3.3 PRACTICAL ALGORITHM

We introduce several improvements over Eq. 9 and propose a practical objective of NFT:

$$\mathcal{L}_D^{\text{NFT}}(\theta) = - \sum_{\mathbf{q}, \mathbf{a}, r} \omega(\mathbf{q}) \sum_t \left[ r \log R_\theta^t(\mathbf{q}, \mathbf{a}) + (1 - r) \log \max_v \left( \frac{1 - \hat{r}_q R_\theta^t(\mathbf{q}, \mathbf{a})}{1 - \hat{r}_q}, \epsilon \right) \right] \quad (10)$$

$$\text{where } R_\theta^t(\mathbf{q}, \mathbf{a}) = \frac{\pi_\theta^+(\mathbf{a}_t|\mathbf{q}, \mathbf{a}_{<t})}{\pi_{\text{old}}(\mathbf{a}_t|\mathbf{q}, \mathbf{a}_{<t})}, \quad \text{and} \quad \hat{r}_q = \frac{1}{K} \sum_{\mathbf{a}|\mathbf{q}} r(\mathbf{q}, \mathbf{a}).$$

Pseudo code is in Algorithm 1. Below, we explain our key design choices.

**Token-level loss.** Eq. 9 essentially deals with sequence data, where answer likelihood  $\pi(\mathbf{a}|\mathbf{q}) = \prod_t \pi(\mathbf{a}_t|\mathbf{q}, \mathbf{a}_{<t})$  is heavily correlated with answer length. This introduces high variance in gradient estimation and causes numerical instabilities during training. Following Schulman et al. (2017); Yu et al. (2025), we view each token decision as an individual unit and sum up their loss in Eq. 10.

**Clipping negative likelihood ratio.** The negative loss in Eq. 10 involves a logarithm whose argument must remain positive, imposing  $(1 - \hat{r}_q R_\theta^t(\mathbf{q}, \mathbf{a})) / (1 - \hat{r}_q) > 0$ . When  $R_\theta^t$  is unoptimized, this requirement may not be satisfied, potentially leading to training collapse. We therefore enforce a minimum value of  $\epsilon > 0$  for the negative likelihood ratio. To preserve gradient flow after clipping, we further apply straight-through gradient (Bengio et al., 2013). Details are in Algorithm 1.

**Prompt weighting.** To focus training on more informative instances, we weight each prompt  $\mathbf{q}$  by  $\omega(\mathbf{q})$ , assigning higher importance to hard questions with low  $r_{\mathbf{q}}$ . An ablation study is posted in Sec. 5.4. This design also helps align NFT with RL algorithms like GRPO (Sec. 4).

#### 4 UNDERSTANDING THE GAP BETWEEN NFT AND GRPO

Despite originating from entirely different theoretical foundations, NFT and GRPO exhibit significant similarities. Notably, we find **GRPO and NFT are equivalent in on-policy training**. To understand this, we calculate and compare their loss gradients:

**Proposition 4.1 (Algorithm Gradient Comparision).** Suppose there are  $\hat{r}_{\mathbf{q}}K$  positive answers and  $(1 - \hat{r}_{\mathbf{q}})K$  negative ones for a given question  $\mathbf{q}$

(a) **GRPO Gradient:** Consider only binary reward  $\{0, 1\}$  in Eq. 3, GRPO loss gradient satisfies

$$\nabla_{\theta} \mathcal{L}_{\mathcal{D}}^{\text{GRPO}}(\theta) = - \sum \left\{ r A_{\mathbf{q}}^{+} \cdot \mathcal{I}[R_{\theta}^t(\mathbf{q}, \mathbf{a}) < 1 + \epsilon'] + (1-r) A_{\mathbf{q}}^{-} \cdot \mathcal{I}[R_{\theta}^t(\mathbf{q}, \mathbf{a}) > 1 - \epsilon'] \right\} \nabla_{\theta} R_{\theta}^t(\mathbf{q}, \mathbf{a}),$$

where  $A_{\mathbf{q}}^{+} = \sqrt{\frac{1-\hat{r}_{\mathbf{q}}}{\hat{r}_{\mathbf{q}}}}$  and  $A_{\mathbf{q}}^{-} = -\sqrt{\frac{\hat{r}_{\mathbf{q}}}{1-\hat{r}_{\mathbf{q}}}}$  are respectively normalized advantages for answers.

(b) **NFT Gradient:** Let  $\omega(\mathbf{q}) = \sqrt{(1 - \hat{r}_{\mathbf{q}})/\hat{r}_{\mathbf{q}}}$ , NFT loss gradient satisfies

$$\nabla_{\theta} \mathcal{L}_{\mathcal{D}}^{\text{NFT}}(\theta) = - \sum \left\{ r A_{\mathbf{q}}^{+} \cdot \frac{1}{R_{\theta}^t(\mathbf{q}, \mathbf{a})} + (1-r) A_{\mathbf{q}}^{-} \cdot \max \left[ \frac{1 - \hat{r}_{\mathbf{q}} R_{\theta}^t(\mathbf{q}, \mathbf{a})}{1 - \hat{r}_{\mathbf{q}}}, \epsilon \right]^{-1} \right\} \nabla_{\theta} R_{\theta}^t(\mathbf{q}, \mathbf{a}).$$

Proofs are provided in Appendix A. Comparing  $\nabla_{\theta} \mathcal{L}_{\mathcal{D}}^{\text{GRPO}}(\theta)$  and  $\nabla_{\theta} \mathcal{L}_{\mathcal{D}}^{\text{NFT}}(\theta)$ , the only difference between GRPO and NFT is their strategy for clipping model gradients when training data becomes *off-policy* (Figure 4). GRPO simply zeros out the gradient when the learned policy  $\pi_{\theta}$  shifts far away from the old policy  $\pi$ , while NFT takes a “softer” decay schedule.

Surprisingly, we find GRPO and NFT to be equivalent when training is totally on-policy, despite their distinctively different derivations:

**Proposition 4.2 (On-policy Gradient Equivalence).** Following the set up of Proposition 4.1 and let  $\epsilon \leq 1$ , GRPO and NFT loss gradient are equivalent in on-policy training:

$$R_{\theta}^t(\mathbf{q}, \mathbf{a}) = 1 \implies \nabla_{\theta} \mathcal{L}_{\mathcal{D}}^{\text{NFT}}(\theta) = \nabla_{\theta} \mathcal{L}_{\mathcal{D}}^{\text{GRPO}}(\theta)$$

**Implicit group normalization.** Proposition 4.1 shows the “normalized advantage” term is implicitly present within NFT’s loss function. This partially justifies the Group Normalization design choice for GRPO, which was initially introduced only as an empirical technique (Shao et al., 2024). In Appendix A, we further demonstrate that by adjusting  $\omega(\mathbf{q}) = 1 - \hat{r}_{\mathbf{q}}$ , NFT also aligns with Dr. GRPO (Liu et al., 2025b). Our findings suggest a strong connection between RL and SL frameworks.

#### 5 EXPERIMENTS

We seek to answer the following questions through our experiments.

1. How does NFT perform in comparison with existing RL algorithms such as GRPO? (Sec. 5.2)
2. How does negative data contribute to NFT’s performance gain? (Sec. 5.3)
3. Which empirical design choices are important to the effectiveness of NFT? (Sec. 5.4)

##### 5.1 EXPERIMENT SETUPS

**Training.** We perform online fine-tuning on Qwen2.5-Math-7B and Qwen2.5-32B (Yang et al., 2024) to enhance their math abilities without relying on external teachers. The training dataset

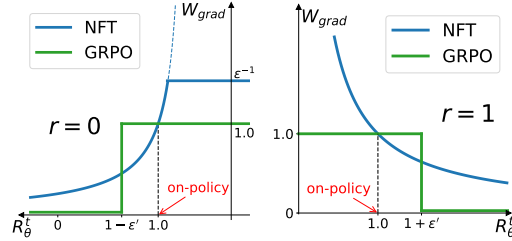


Figure 4: Gradient weight for NFT and GRPO.

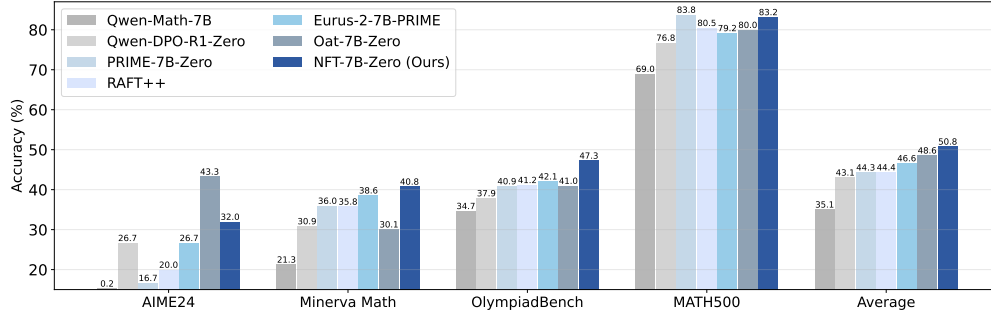


Figure 5: Comparison of the released NFT-7B with other zero-style math models of Qwen series.

Table 1: NFT performs competitively compared with other algorithms. We report avg@32 for AIME24, AIME25, and AMC23 and avg@1 for others. Numbers within 1 % of the max are bolded.

Model	AIME24	MATH500	AIME25	AMC23	Olympiad	Minerva	Average
Qwen2.5-Math-7B	13.3	69.0	5.5	45.8	34.7	21.3	31.6
<i>Preference fine-tuning</i>							
+ DPO	29.8	79.8	13.8	83.2	48.0	39.0	48.9
<i>Reinforcement fine-tuning</i>							
+ GRPO	30.2	80.4	17.1	79.5	<b>51.8</b>	38.2	49.5
+ Dr. GRPO	31.8	<b>83.4</b>	15.7	80.2	49.6	38.2	49.8
+ DAPO	<b>33.1</b>	81.6	<b>18.7</b>	85.0	49.9	39.3	<b>51.2</b>
<i>Supervised fine-tuning</i>							
+ RFT	<b>33.7</b>	79.8	13.4	79.7	44.3	38.6	48.3
+ NFT	32.0	<b>83.2</b>	<b>18.3</b>	<b>88.5</b>	47.3	<b>40.8</b>	<b>51.7</b>
Qwen2.5-72B	4.1	68.6	1.0	45.0	31.1	27.9	29.6
+ DAPO	<b>44.1</b>	<b>89.2</b>	<b>33.4</b>	90.9	<b>54.1</b>	47.5	<b>59.9</b>
+ RFT	29.9	86.2	19.1	92.4	45.3	44.1	52.8
+ NFT	37.8	<b>88.4</b>	31.5	<b>93.8</b>	<b>55.0</b>	<b>48.9</b>	<b>59.2</b>

is the publicly available DAPO-Math-17k (Yu et al., 2025), which consists solely of math questions paired with ground-truth answers in integer form. During training, all models are fine-tuned for approximately 5,000 gradient steps with a batch size of 512. Generation temperature is 1.0.

**Evaluation.** We evaluate models on six validation benchmarks and report their average accuracy: AIME 2024, AIME 2025, AMC 2023 (Li et al., 2024), MATH500 (Hendrycks et al., 2021), OlympiadBench (He et al., 2024), and Minerva Math (Lewkowycz et al., 2022).

**Baseline methods.** We compare against a set of online fine-tuning algorithms, including Iterative DPO (Rafailov et al., 2023; Xiong et al., 2023; Guo et al., 2024), GRPO (Shao et al., 2024), Dr. GRPO (Liu et al., 2025b), DAPO (Yu et al., 2025), and RFT (Dong et al., 2023; Yuan et al., 2023b).

## 5.2 NFT PERFORMANCE EVALUATION

**Model comparison.** By applying NFT to Qwen2.5-Math-7B, we release NFT-7B-Zero (Figure 5). NFT-7B-Zero achieves competitive performance on all benchmarks compared to other zero-style 7B math models (Cui et al., 2025; Liu et al., 2025b; Xiong et al., 2025; 2023). This provides strong empirical evidence for the effectiveness of the NFT algorithm and demonstrates that SL alone can enable effective verification-driven training in math tasks.

**Algorithm comparison.** To isolate the contribution of the algorithm itself, we benchmarked various online algorithms using identical training data, infrastructure, and general hyperparameters (Table 1). Results show that NFT matches or even surpasses state-of-the-art methods such as DAPO. Figure 6 and 11 present training curves across multiple runs. NFT exhibits convergence speed and final performance on par with DAPO, further supporting its stability.

## 5.3 BENEFITS OF NEGATIVE DATA



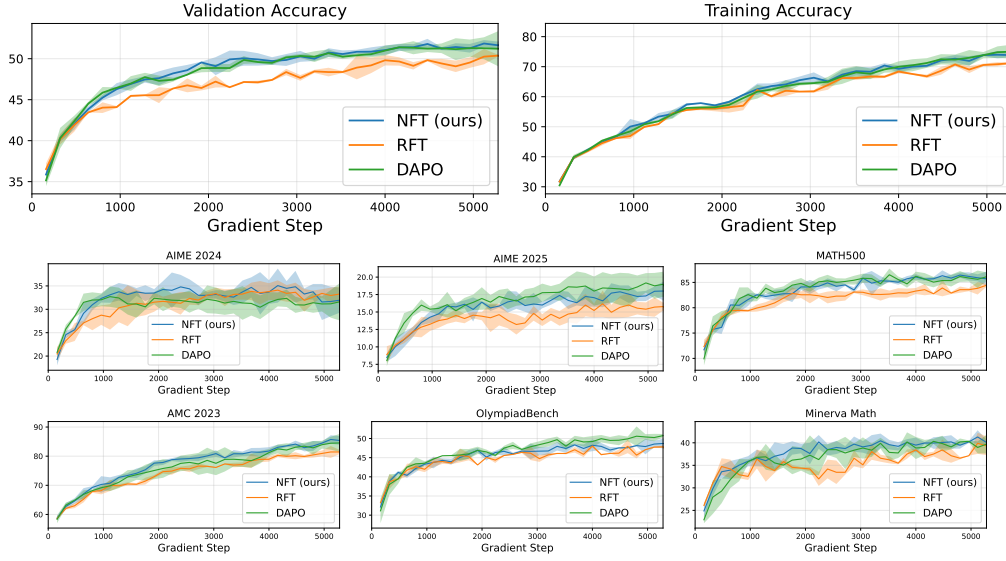


Figure 6: Training and validation accuracy curves for 7B experiments. We conducted 3-4 random and independent experiments for each algorithm and report their mean  $\pm$  std results.

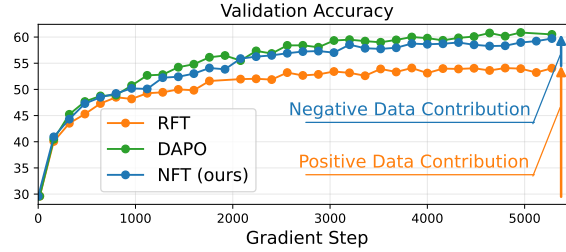


Figure 7: Average accuracy across 6 benchmarks for 32B experiments. More curves in Appendix D.

### Negative feedback enhances performance and exploration.

Table 1 shows that NFT consistently outperforms RFT by a clear margin, highlighting the benefit of incorporating negative feedback during training. Notably, we observe a clear divergence in training dynamics between RFT and NFT. Across both 7B and 32B settings, RFT tends to reduce entropy over time, whereas NFT and RL methods like DAPO encourage increasing entropy (Figure 8). This suggests more exploration (Yu et al., 2025), potentially hindering why NFT outperforms RFT.

### Negative feedback becomes increasingly important in larger models.

The performance gap between RFT and NFT widens faster over training in 32B experiments compared with the trend in 7B (Figure 11). Similarly, DeepSeek-AI (2025) also notes RL offers greater benefits over SFT in larger models. A potential explanation could be the increasing importance of negative data: Larger models already memorize well enough, so the ability to reflect on mistakes becomes a new bottleneck.

### RFT remains a strong baseline.

Although surpassed by numerous algorithms, RFT still deserves attention due to its extreme simplicity. In 32B settings (Figure 11), learning from positive data (RFT) contributes to 80% of the total gain achieved by our best-performing model, while negative data only accounts for the remaining 20%. These findings echo recent studies (Yue et al., 2025; Xiong et al., 2025; Liu et al., 2025a; Zhao et al., 2025; Wang et al., 2025), which suggest RL primarily amplifies existing capabilities in large models rather than fostering new skills. How to exploit negative feedback remains an open challenge with heavy potential.

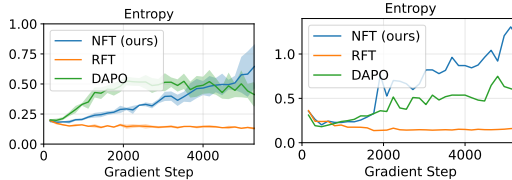


Figure 8: Entropy curves for 7B and 32B runs.



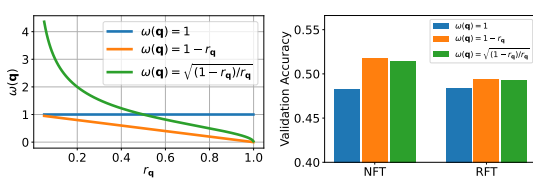
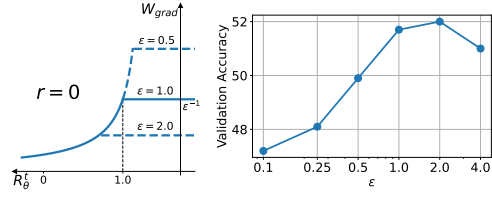


Figure 9: Effect of prompt weighting.

Figure 10: Effect of negative ratio clip value  $\epsilon$ .

#### 5.4 INGREDIENTS BEHIND NFT’S EFFECTIVENESS

We discuss two empirical design choices that notably help NFT achieve strong performance.

**Prioritize harder questions.** We find that assigning higher weights to difficult questions with a low answer correctness rate  $\hat{r}_q$  can enhance model performance. We achieve this mainly by selecting  $\omega(q)$  in Eq. 10 from 3 choices: (1)  $\omega(q) = 1$ . (2)  $\omega(q) = 1 - r_q$ , which aligns NFT with Dr. GRPO in on-policy training (Sec. 4). (3)  $\omega(q) = \sqrt{(1 - r_q)/r_q}$ , which aligns NFT with GRPO. Figure 9 visualizes different  $\omega(q)$ , and their effect for NFT and RFT. We find choices (2) and (3) perform similarly, both outperforming constant weighting choice (1).

**Avoid overpenalizing mistakes.** The clip value  $\epsilon$  of NFT (Eq. 10) sets an upper bound on the penalty weight when the likelihood ratio  $R_\theta^t$  for negative answers increases. When  $\epsilon$  is small (e.g., near zero), the algorithm empirically assigns high penalties to rising likelihoods of incorrect answers (Figure 10). However, our experiments show that overly aggressive penalization with  $\epsilon \rightarrow 0$  degrades overall performance. We thus adopt a default setting of  $\epsilon = 1.0$ .

## 6 RELATED WORKS

Reinforcement Learning with Verifiable Rewards (RLVR) has advanced the frontier of LLM reasoning (DeepSeek-AI, 2025; OpenAI, 2024; Team et al., 2025; Chen et al., 2025c). Compared with previous RL practices that rely on strong reward models (Wang et al., 2023b; Yuan et al., 2024; Zhang et al., 2024) to simulate human feedback (Ouyang et al., 2022; Christiano et al., 2017; Cui et al., 2025), RLVR turns to a ground truth verifier for providing reliable binary supervision (Lambert et al., 2024; Shao et al., 2024). Moreover, unlike preference-based learning algorithms such as DPO (Rafailov et al., 2023; Cai et al., 2023; Azar et al., 2024; Ethayarajh et al., 2024; Wang et al., 2023a; Chen et al., 2024; Hong et al., 2024; Xu et al., 2023), RLVR does not require paired preference data, rendering it more flexible and memory-efficient. Despite the demonstrated effectiveness of RL algorithms in verification-driven training (Li et al., 2023; Ahmadian et al., 2024; Hu, 2025; Yu et al., 2025; Chu et al., 2025b; Yuan et al., 2025; Yan et al., 2025), recent studies suggest that supervised learning (SL) may also suffice for achieving verification-driven training in LLMs (Dong et al., 2023). Our method further addresses SL’s inability to incorporate negative feedback (Hua et al., 2024), bridging both the theoretical and the performance gap between the two fields.

A key design of NFT involves implicit policy modeling for direct policy optimization. This design, emphasizing direct optimization via implicitly defined models, shares conceptual similarities with some existing approaches. In preference-based training, DPO (Rafailov et al., 2023) introduces an implicit reward model parameterized by the policy network to allow optimizing policies directly. Recent visual modeling efforts also leverage implicit conditional or residual models parameterized by generation networks to avoid guided sampling (Chen et al., 2025b;a; Zheng et al., 2025).

## 7 CONCLUSION

In this work, we introduce Negative-aware Fine-Tuning (NFT), a supervised approach that enables LLMs to learn from their own negative generations. In online training, NFT substantially improves upon supervised learning baselines through the additional leverage of negative feedback, achieving performance comparable to leading RL algorithms like GRPO. Notably, we unveiled a theoretical equivalence between NFT and GRPO under strict-on-policy conditions. These findings bridge the conceptual and practical gap between SL and RL paradigms in verification-driven training.

## REPRODUCIBILITY

We submit code in the supplementary material. Code, model weights, and dataset will be released.

## THE USE OF LARGE LANGUAGE MODELS (LLMs)

We used large language models (LLMs) solely as a writing assistant for language polishing and improving clarity of presentation. The LLMs were not involved in research ideation, methodological design, experimental execution, or result analysis. All scientific contributions and substantive writing were carried out by the authors.

## REFERENCES

- Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv preprint arXiv:2402.14740*, 2024.
- Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from human preferences. In *AISTATS*, 2024.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- Tianchi Cai, Xierui Song, Jiyan Jiang, Fei Teng, Jinjie Gu, and Guannan Zhang. Ulma: Unified language model alignment with demonstration and point-wise human preference. *arXiv preprint arXiv:2312.02554*, 2023.
- Huayu Chen, Guande He, Lifan Yuan, Ganqu Cui, Hang Su, and Jun Zhu. Noise contrastive alignment of language models with explicit rewards. In *NeurIPS*, 2024.
- Huayu Chen, Kai Jiang, Kaiwen Zheng, Jianfei Chen, Hang Su, and Jun Zhu. Visual generation without guidance. In *ICML*, 2025a.
- Huayu Chen, Hang Su, Peize Sun, and Jun Zhu. Toward guidance-free ar visual generation via condition contrastive alignment. In *ICLR*, 2025b.
- Yang Chen, Zhuolin Yang, Zihan Liu, Chankyu Lee, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. Acereason-nemotron: Advancing math and code reasoning through reinforcement learning. *arXiv preprint*, 2025c.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In *NeurIPS*, 2017.
- 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*, 2025a.
- Xiangxiang Chu, Hailang Huang, Xiao Zhang, Fei Wei, and Yong Wang. Gpg: A simple and strong reinforcement learning baseline for model reasoning. *arXiv preprint arXiv:2504.02546*, 2025b.
- Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang, Wendi Li, Bingxiang He, Yuchen Fan, Tianyu Yu, Qixin Xu, Weize Chen, et al. Process reinforcement through implicit rewards. *arXiv preprint arXiv:2502.01456*, 2025.
- DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

- Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *TMLR*, 2023.
- Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexandre Rame, Thomas Mesnard, Yao Zhao, Bilal Piot, et al. Direct language model alignment from online ai feedback. *arXiv preprint arXiv:2402.04792*, 2024.
- 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, 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.
- Jiwoo Hong, Noah Lee, and James Thorne. Orpo: Monolithic preference optimization without reference model. *arXiv preprint arXiv:2403.07691*, 2024.
- Jian Hu. Reinforce++: A simple and efficient approach for aligning large language models. *arXiv preprint arXiv:2501.03262*, 2025.
- Ermo Hua, Bqing Qi, Kaiyan Zhang, Yue Yu, Ning Ding, Xingtai Lv, Kai Tian, and Bowen Zhou. Intuitive fine-tuning: Towards simplifying alignment into a single process. *arXiv preprint arXiv:2405.11870*, 2024.
- Hynek Kydlíček. Math-verify: Math verification library, 2024. URL <https://github.com/huggingface/math-verify>.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.
- Sergey Levine. Reinforcement learning and control as probabilistic inference: Tutorial and review. *arXiv preprint arXiv:1805.00909*, 2018.
- 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. In *NeurIPS*, 2022.
- Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Huang, Kashif Rasul, Longhui Yu, Albert Q Jiang, Ziju Shen, et al. Numinamath: The largest public dataset in ai4maths with 860k pairs of competition math problems and solutions. *Hugging Face repository*, 2024.
- Ziniu Li, Tian Xu, Yushun Zhang, Zhihang Lin, Yang Yu, Ruoyu Sun, and Zhi-Quan Luo. Remax: A simple, effective, and efficient reinforcement learning method for aligning large language models. *arXiv preprint arXiv:2310.10505*, 2023.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- Zichen Liu, Changyu Chen, Wenjun Li, Tianyu Pang, Chao Du, and Min Lin. There may not be aha moment in r1-zero-like training—a pilot study, 2025a.
- Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding r1-zero-like training: A critical perspective. *arXiv preprint arXiv:2503.20783*, 2025b.

- NVIDIA. Cosmos-reason1: From physical common sense to embodied reasoning. *arXiv preprint arXiv:2503.15558*, 2025.
- OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- OpenAI. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022.
- 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. In *NeurIPS*, 2023.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *ICML*, 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, Mingchuan Zhang, YK Li, Y Wu, and Daya Guo. 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.
- Feifan Song, Bowen Yu, Minghao Li, Haiyang Yu, Fei Huang, Yongbin Li, and Houfeng Wang. Preference ranking optimization for human alignment. *arXiv preprint arXiv:2306.17492*, 2023.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *NeurIPS*, 1999.
- Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Chaoqi Wang, Yibo Jiang, Chenghao Yang, Han Liu, and Yuxin Chen. Beyond reverse kl: Generalizing direct preference optimization with diverse divergence constraints. *arXiv preprint arXiv:2309.16240*, 2023a.
- Peiyi Wang, Lei Li, Zhihong Shao, RX Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations. *arXiv preprint arXiv:2312.08935*, 2023b.
- Yiping Wang, Qing Yang, Zhiyuan Zeng, Liliang Ren, Lucas Liu, Baolin Peng, Hao Cheng, Xuehai He, Kuan Wang, Jianfeng Gao, et al. Reinforcement learning for reasoning in large language models with one training example. *arXiv preprint arXiv:2504.20571*, 2025.
- Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. Iterative preference learning from human feedback: Bridging theory and practice for rlhf under kl-constraint. *arXiv preprint arXiv:2312.11456*, 2023.
- Wei Xiong, Jiarui Yao, Yuhui Xu, Bo Pang, Lei Wang, Doyen Sahoo, Junnan Li, Nan Jiang, Tong Zhang, Caiming Xiong, et al. A minimalist approach to llm reasoning: from rejection sampling to reinforce. *arXiv preprint arXiv:2504.11343*, 2025.

- Jing Xu, Andrew Lee, Sainbayar Sukhbaatar, and Jason Weston. Some things are more cringe than others: Iterative preference optimization with the pairwise cringe loss. *arXiv preprint arXiv:2312.16682*, 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, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2. 5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, et al. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*, 2025.
- Hongyi Yuan, Zheng Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. Rrhf: Rank responses to align language models with human feedback. In *NeurIPS*, 2023a.
- Lifan Yuan, Wendi Li, Huayu Chen, Ganqu Cui, Ning Ding, Kaiyan Zhang, Bowen Zhou, Zhiyuan Liu, and Hao Peng. Free process rewards without process labels. *arXiv preprint arXiv:2412.01981*, 2024.
- Yurun Yuan, Fan Chen, Zeyu Jia, Alexander Rakhlin, and Tengyang Xie. Trajectory bellman residual minimization: A simple value-based method for llm reasoning, 2025. URL <https://arxiv.org/abs/2505.15311>.
- Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Chuanqi Tan, and Chang Zhou. Scaling relationship on learning mathematical reasoning with large language models. *arXiv preprint arXiv:2308.01825*, 2023b.
- Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Shiji Song, and Gao Huang. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model? *arXiv preprint arXiv:2504.13837*, 2025.
- Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the wild. *arXiv preprint arXiv:2503.18892*, 2025.
- Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. Generative verifiers: Reward modeling as next-token prediction. *arXiv preprint arXiv:2408.15240*, 2024.
- Rosie Zhao, Alexandru Meterez, Sham Kakade, Cengiz Pehlevan, Samy Jelassi, and Eran Malach. Echo chamber: RL post-training amplifies behaviors learned in pretraining. *arXiv preprint arXiv:2504.07912*, 2025.
- Kaiwen Zheng, Yongxin Chen, Huayu Chen, Guande He, Ming-Yu Liu, Jun Zhu, and Qinsheng Zhang. Direct discriminative optimization: Your likelihood-based visual generative model is secretly a gan discriminator. In *ICML*, 2025.

## A PROOF OF THEOREMS

**Theorem A.1 (Policy Optimization with Negative Answers).** Consider the maximum-likelihood objective for training the implicit negative policy  $\pi_{\theta}^{-}$ :

$$\max_{\theta} \mathbb{E}_{p(\mathbf{q})\pi^{-}(\mathbf{a}|\mathbf{q})} [\log \pi_{\theta}^{-}(\mathbf{a}|\mathbf{q})] \Leftrightarrow \min_{\theta} \left[ -\mathbb{E}_{(\mathbf{q},\mathbf{a}) \sim \mathcal{D}^{-}} \log \frac{\pi_{\text{old}}(\mathbf{a}|\mathbf{q}) - r_{\mathbf{q}}\pi_{\theta}^{+}(\mathbf{a}|\mathbf{q})}{1 - r_{\mathbf{q}}} \right]$$

Assuming unlimited data and model capacity, the optimal solution for solving Eq. 8 is

$$\forall \mathbf{q}, \mathbf{a} : \quad \pi_{\theta}^{+}(\mathbf{a}|\mathbf{q}) = \pi^{+}(\mathbf{a}|\mathbf{q})$$

*Proof.* The proof is quite straightforward. First, we show that maximum-likelihood training leads to the optimal solution  $\pi_{\theta^{*}}^{-}(\mathbf{a}|\mathbf{q}) = \pi^{-}(\mathbf{a}|\mathbf{q})$ .

$$\begin{aligned} & \max_{\theta} \mathbb{E}_{p(\mathbf{q})\pi^{-}(\mathbf{a}|\mathbf{q})} [\log \pi_{\theta}^{-}(\mathbf{a}|\mathbf{q})] \\ & \Leftrightarrow \max_{\theta} \mathbb{E}_{p(\mathbf{q})\pi^{-}(\mathbf{a}|\mathbf{q})} [\log \pi_{\theta}^{-}(\mathbf{a}|\mathbf{q}) - \log \pi^{-}(\mathbf{a}|\mathbf{q})] \\ & \Leftrightarrow \min_{\theta} \mathbb{E}_{p(\mathbf{q})} D_{\text{KL}} [\pi^{-}(\mathbf{a}|\mathbf{q}) \| \pi_{\theta}^{-}(\mathbf{a}|\mathbf{q})] \end{aligned}$$

Since  $D_{\text{KL}} [\pi^{-}(\mathbf{a}|\mathbf{q}) \| \pi_{\theta}^{-}(\mathbf{a}|\mathbf{q})] \geq 0$ . The equality holds if and only if  $\forall \mathbf{q} : \pi_{\theta}^{-} = \pi^{-}$ . We thus have

$$\forall \mathbf{q}, \mathbf{a} : \quad \pi_{\theta^{*}}^{-}(\mathbf{a}|\mathbf{q}) = \pi^{-}(\mathbf{a}|\mathbf{q}). \quad (11)$$

Next, we prove  $\pi_{\theta^{*}}^{+} = \pi^{+}$ .

Note that that during training,  $\pi_{\theta}^{-}$  is only an *implicit* policy defined by  $\pi_{\theta}^{+}$  through

$$\pi_{\theta}^{-}(\mathbf{a}|\mathbf{q}) := \frac{\pi_{\text{old}}(\mathbf{a}|\mathbf{q}) - r_{\mathbf{q}}\pi_{\theta}^{+}(\mathbf{a}|\mathbf{q})}{1 - r_{\mathbf{q}}}.$$

On the other hand, the negative data distribution  $\pi^{-}$  has the same relationship with  $\pi^{+}$  by Eq. 7.

$$\pi^{-}(\mathbf{a}|\mathbf{q}) := \frac{\pi_{\text{old}}(\mathbf{a}|\mathbf{q}) - r_{\mathbf{q}}\pi^{+}(\mathbf{a}|\mathbf{q})}{1 - r_{\mathbf{q}}}.$$

We have ensured  $0 < r_{\mathbf{q}} < 1$  during training, combining these observations and Eq. 11, we have

$$\forall \mathbf{q}, \mathbf{a} : \quad \pi_{\theta^{*}}^{+}(\mathbf{a}|\mathbf{q}) = \pi^{+}(\mathbf{a}|\mathbf{q})$$

□

**Proposition A.2 (Algorithm Gradient Comparision).** Suppose there are  $\hat{r}_{\mathbf{q}}K$  positive answers and  $(1 - \hat{r}_{\mathbf{q}})K$  negative ones for a given question  $\mathbf{q}$

**(a) GRPO Gradient:** Consider only binary reward  $\{0, 1\}$  in Eq. 3, GRPO loss gradient satisfies

$$\nabla_{\theta} \mathcal{L}_{\mathcal{D}}^{\text{GRPO}}(\theta) = - \sum \left\{ r A_{\mathbf{q}}^{+} \cdot \mathcal{I} [R_{\theta}^t(\mathbf{q}, \mathbf{a}) < 1 + \epsilon'] + (1-r) A_{\mathbf{q}}^{-} \cdot \mathcal{I} [R_{\theta}^t(\mathbf{q}, \mathbf{a}) > 1 - \epsilon'] \right\} \nabla_{\theta} R_{\theta}^t(\mathbf{q}, \mathbf{a}),$$

where  $A_{\mathbf{q}}^{+} = \sqrt{\frac{1-\hat{r}_{\mathbf{q}}}{\hat{r}_{\mathbf{q}}}}$  and  $A_{\mathbf{q}}^{-} = -\sqrt{\frac{\hat{r}_{\mathbf{q}}}{1-\hat{r}_{\mathbf{q}}}}$  are respectively normalized advantages for answers.

**(b) NFT Gradient:** Let  $\omega(\mathbf{q}) = \sqrt{(1 - \hat{r}_{\mathbf{q}})/\hat{r}_{\mathbf{q}}}$ , NFT loss gradient satisfies

$$\nabla_{\theta} \mathcal{L}_{\mathcal{D}}^{\text{NFT}}(\theta) = - \sum \left\{ r A_{\mathbf{q}}^{+} \cdot \frac{1}{R_{\theta}^t(\mathbf{q}, \mathbf{a})} + (1-r) A_{\mathbf{q}}^{-} \cdot \max \left[ \frac{1 - \hat{r}_{\mathbf{q}} R_{\theta}^t(\mathbf{q}, \mathbf{a})}{1 - \hat{r}_{\mathbf{q}}}, \epsilon \right]^{-1} \right\} \nabla_{\theta} R_{\theta}^t(\mathbf{q}, \mathbf{a}).$$

*Proof.* **(a) GRPO Gradient:** We first copy down the GRPO loss definition from Eq. 3.

$$\mathcal{L}_{\mathcal{D}}^{\text{GRPO}}(\theta) = - \sum_{\mathbf{q}, \mathbf{a}, t} \min \left[ R_{\theta}^t(\mathbf{q}, \mathbf{a}) \hat{A}_{\mathbf{q}, \mathbf{a}}, \text{clip}(R_{\theta}^t(\mathbf{q}, \mathbf{a}), 1 - \epsilon', 1 + \epsilon') \hat{A}_{\mathbf{q}, \mathbf{a}} \right].$$

The normalized advantage can be estimated as

$$\begin{aligned}\hat{A}_{\mathbf{q},\mathbf{a}} &:= [r(\mathbf{q},\mathbf{a}) - \text{mean}\{r^{1:K}\}] / \text{std}\{r^{1:K}\} \\ \text{where } \text{mean}\{r^{1:K}\} &= \frac{1}{K} [\hat{r}_{\mathbf{q}}K * 1 + (1 - \hat{r}_{\mathbf{q}})K * 0] = \hat{r}_{\mathbf{q}} \\ \text{and } \text{std}\{r^{1:K}\} &= \sqrt{\frac{1}{K} [\hat{r}_{\mathbf{q}}K * (1 - \hat{r}_{\mathbf{q}})^2 + (1 - \hat{r}_{\mathbf{q}})K * (0 - \hat{r}_{\mathbf{q}})^2]} = \sqrt{\hat{r}_{\mathbf{q}}(1 - \hat{r}_{\mathbf{q}})}.\end{aligned}$$

When  $\mathbf{a}$  is a positive answer,  $r(\mathbf{q},\mathbf{a}) = 1$ . We have  $A_{\mathbf{q}}^+ = \sqrt{\frac{1 - \hat{r}_{\mathbf{q}}}{\hat{r}_{\mathbf{q}}}} > 0$ .

$$\begin{aligned}\mathcal{L}_{\mathcal{D}^+}^{\text{GRPO}}(\theta) &= - \sum_{\mathbf{q},\mathbf{a}^+,t} \min [R_{\theta}^t(\mathbf{q},\mathbf{a}^+), 1 + \epsilon'] \hat{A}_{\mathbf{q},\mathbf{a}^+} \\ \nabla_{\theta} \mathcal{L}_{\mathcal{D}^+}^{\text{GRPO}}(\theta) &= - \sum_{\mathbf{q},\mathbf{a}^+,t} \mathbf{A}_{\mathbf{q}}^+ \cdot \mathcal{I} [R_{\theta}^t(\mathbf{q},\mathbf{a}^+) < 1 + \epsilon'].\end{aligned}\quad (12)$$

When  $\mathbf{a}$  is a negative answer,  $r(\mathbf{q},\mathbf{a}) = 0$ . We have  $A_{\mathbf{q}}^- = -\sqrt{\frac{\hat{r}_{\mathbf{q}}}{1 - \hat{r}_{\mathbf{q}}}} < 0$ .

$$\begin{aligned}\mathcal{L}_{\mathcal{D}^-}^{\text{GRPO}}(\theta) &= - \sum_{\mathbf{q},\mathbf{a}^-,t} \max [R_{\theta}^t(\mathbf{q},\mathbf{a}^-), 1 - \epsilon'] \hat{A}_{\mathbf{q},\mathbf{a}^-} \\ \nabla_{\theta} \mathcal{L}_{\mathcal{D}^-}^{\text{GRPO}}(\theta) &= - \sum_{\mathbf{q},\mathbf{a}^-,t} \mathbf{A}_{\mathbf{q}}^- \cdot \mathcal{I} [R_{\theta}^t(\mathbf{q},\mathbf{a}^-) > 1 - \epsilon'].\end{aligned}\quad (13)$$

Combining Eq. 12 and Eq. 13, we have

$$\nabla_{\theta} \mathcal{L}_{\mathcal{D}}^{\text{GRPO}}(\theta) = - \sum \left\{ r \mathbf{A}_{\mathbf{q}}^+ \cdot \mathcal{I} [R_{\theta}^t(\mathbf{q},\mathbf{a}) < 1 + \epsilon'] + (1-r) \mathbf{A}_{\mathbf{q}}^- \cdot \mathcal{I} [R_{\theta}^t(\mathbf{q},\mathbf{a}) > 1 - \epsilon'] \right\} \nabla_{\theta} R_{\theta}^t(\mathbf{q},\mathbf{a}),$$

**(a) NFT Gradient:** We copy down the NFT loss definition from Eq. 10.

$$\mathcal{L}_{\mathcal{D}}^{\text{NFT}}(\theta) = - \sum_{\mathbf{q},\mathbf{a},t} \omega(\mathbf{q}) \left[ r \log R_{\theta}^t(\mathbf{q},\mathbf{a}) + (1-r) \log \max_{\vee} \left( \frac{1 - \hat{r}_{\mathbf{q}} R_{\theta}^t(\mathbf{q},\mathbf{a})}{1 - \hat{r}_{\mathbf{q}}}, \epsilon \right) \right]$$

When  $\mathbf{a}$  is a positive answer,  $r(\mathbf{q},\mathbf{a}) = 1$ .

$$\begin{aligned}\mathcal{L}_{\mathcal{D}^+}^{\text{NFT}}(\theta) &= - \sum_{\mathbf{q},\mathbf{a}^+,t} \omega(\mathbf{q}) \log R_{\theta}^t(\mathbf{q},\mathbf{a}) \\ &= - \sum_{\mathbf{q},\mathbf{a}^+,t} \sqrt{\frac{1 - \hat{r}_{\mathbf{q}}}{\hat{r}_{\mathbf{q}}}} \log R_{\theta}^t(\mathbf{q},\mathbf{a}) \\ &= - \sum_{\mathbf{q},\mathbf{a}^+,t} \mathbf{A}_{\mathbf{q}}^+ \log R_{\theta}^t(\mathbf{q},\mathbf{a}^+)\end{aligned}$$

$$\nabla_{\theta} \mathcal{L}_{\mathcal{D}^+}^{\text{NFT}} = - \sum_{\mathbf{q},\mathbf{a}^+,t} \mathbf{A}_{\mathbf{q}}^+ \frac{1}{R_{\theta}^t(\mathbf{q},\mathbf{a}^+)} \nabla_{\theta} R_{\theta}^t(\mathbf{q},\mathbf{a}^+)\quad (14)$$

When  $\mathbf{a}$  is a negative answer,  $r(\mathbf{q},\mathbf{a}) = 0$ .

$$\begin{aligned}\mathcal{L}_{\mathcal{D}^-}^{\text{NFT}}(\theta) &= - \sum_{\mathbf{q},\mathbf{a}^-,t} \omega(\mathbf{q}) \log \left[ \max_{\vee} \left( \frac{1 - \hat{r}_{\mathbf{q}} R_{\theta}^t(\mathbf{q},\mathbf{a}^-)}{1 - \hat{r}_{\mathbf{q}}}, \epsilon \right) \right] \\ &= - \sum_{\mathbf{q},\mathbf{a}^-,t} \sqrt{\frac{1 - \hat{r}_{\mathbf{q}}}{\hat{r}_{\mathbf{q}}}} \log \left[ \max_{\vee} \left( \frac{1 - \hat{r}_{\mathbf{q}} R_{\theta}^t(\mathbf{q},\mathbf{a}^-)}{1 - \hat{r}_{\mathbf{q}}}, \epsilon \right) \right]\end{aligned}$$



$$\begin{aligned}
\nabla_{\theta} \mathcal{L}_{\mathcal{D}^{-}}^{\text{NFT}} &= - \sum_{\mathbf{q}, \mathbf{a}^{-}, t} \sqrt{\frac{1 - \hat{r}_{\mathbf{q}}}{\hat{r}_{\mathbf{q}}}} \left[ \max\left(\frac{1 - \hat{r}_{\mathbf{q}} R_{\theta}^t(\mathbf{q}, \mathbf{a}^{-})}{1 - \hat{r}_{\mathbf{q}}}, \epsilon\right)^{-1} \cdot \frac{-\hat{r}_{\mathbf{q}}}{1 - \hat{r}_{\mathbf{q}}} \cdot \nabla_{\theta} R_{\theta}^t(\mathbf{q}, \mathbf{a}^{-}) \right] \\
&= - \sum_{\mathbf{q}, \mathbf{a}^{-}, t} - \sqrt{\frac{\hat{r}_{\mathbf{q}}}{1 - \hat{r}_{\mathbf{q}}}} \left[ \max\left(\frac{1 - \hat{r}_{\mathbf{q}} R_{\theta}^t(\mathbf{q}, \mathbf{a}^{-})}{1 - \hat{r}_{\mathbf{q}}}, \epsilon\right)^{-1} \cdot \nabla_{\theta} R_{\theta}^t(\mathbf{q}, \mathbf{a}^{-}) \right] \\
&= - \sum_{\mathbf{q}, \mathbf{a}^{-}, t} A_{\mathbf{q}}^{-} \left[ \max\left(\frac{1 - \hat{r}_{\mathbf{q}} R_{\theta}^t(\mathbf{q}, \mathbf{a}^{-})}{1 - \hat{r}_{\mathbf{q}}}, \epsilon\right)^{-1} \cdot \nabla_{\theta} R_{\theta}^t(\mathbf{q}, \mathbf{a}^{-}) \right] \tag{15}
\end{aligned}$$

Combining Eq. 14 and Eq. 15, we have

$$\nabla_{\theta} \mathcal{L}_{\mathcal{D}}^{\text{NFT}}(\theta) = - \sum \left\{ r A_{\mathbf{q}}^{+} \cdot \frac{1}{R_{\theta}^t(\mathbf{q}, \mathbf{a})} + (1 - r) A_{\mathbf{q}}^{-} \cdot \max\left[\frac{1 - \hat{r}_{\mathbf{q}} R_{\theta}^t(\mathbf{q}, \mathbf{a})}{1 - \hat{r}_{\mathbf{q}}}, \epsilon\right]^{-1} \right\} \nabla_{\theta} R_{\theta}^t(\mathbf{q}, \mathbf{a}).$$

□

**Remark A.3 (Dr. GRPO).** The main practical difference between Dr. GRPO (Liu et al., 2025b) and GRPO is that Dr. GRPO removes the std normalization term when estimating group-normalized advantages. Following Proposition 4.1, we simply need to set  $\omega(\mathbf{q}) = 1 - \hat{r}_{\mathbf{q}}$  instead of  $\omega(\mathbf{q}) = \sqrt{\frac{1 - \hat{r}_{\mathbf{q}}}{\hat{r}_{\mathbf{q}}}}$  to align with the Dr. GRPO loss function.

**Proposition A.4 (On-policy Gradient Equivalence).** *Following the set up of Proposition 4.1 and let  $\epsilon \leq 1$ , GRPO and NFT loss gradient are equivalent in on-policy training:*

$$R_{\theta}^t(\mathbf{q}, \mathbf{a}) = 1 \implies \nabla_{\theta} \mathcal{L}_{\mathcal{D}}^{\text{NFT}}(\theta) = \nabla_{\theta} \mathcal{L}_{\mathcal{D}}^{\text{GRPO}}(\theta)$$

*Proof.* The proof is simple. When on-policy,  $R_{\theta}^t(\mathbf{q}, \mathbf{a}) = 1$ .

Regarding positive answers  $\mathbf{a}^{+}$ , GRPO gradient (Eq. 12) and NFT gradient (Eq. 14) both become

$$\nabla_{\theta} \mathcal{L}_{\mathcal{D}^{+}}^{\text{GRPO}}(\theta) = \nabla_{\theta} \mathcal{L}_{\mathcal{D}^{+}}^{\text{NFT}}(\theta) = A_{\mathbf{q}}^{+} \nabla_{\theta} R_{\theta}^t(\mathbf{q}, \mathbf{a}^{+}).$$

Regarding negative answers  $\mathbf{a}^{-}$ , GRPO gradient (Eq. 13) and NFT gradient (Eq. 15) both become

$$\nabla_{\theta} \mathcal{L}_{\mathcal{D}^{-}}^{\text{GRPO}}(\theta) = \nabla_{\theta} \mathcal{L}_{\mathcal{D}^{-}}^{\text{NFT}}(\theta) = A_{\mathbf{q}}^{-} \nabla_{\theta} R_{\theta}^t(\mathbf{q}, \mathbf{a}^{-}).$$

□

## B DISCUSSION FOR CONTINUOUS REWARDS

In this section, we explain how NFTs can be extended to continuous reward settings.

First, we define  $r := p(\mathbf{o} = 1 | \mathbf{q}, \mathbf{a}) \in [0, 1]$ , where  $\mathbf{o} = 1$  means the answer is an optimal answer.

Then we define positive and negative distributions.

$$\begin{aligned}
\pi^{+}(\mathbf{a} | \mathbf{q}) &:= \pi(\mathbf{a} | \mathbf{q}, \mathbf{o} = 1) = \frac{\pi_{\text{old}}(\mathbf{a} | \mathbf{q}) p(\mathbf{o} = 1 | \mathbf{q}, \mathbf{a})}{\sum_{\mathbf{A}} \pi_{\text{old}}(\mathbf{a} | \mathbf{q}) p(\mathbf{o} = 1 | \mathbf{q}, \mathbf{a})} \propto \pi_{\text{old}}(\mathbf{a} | \mathbf{q}) r(\mathbf{q}, \mathbf{a}), \\
\pi^{-}(\mathbf{a} | \mathbf{q}) &:= \pi(\mathbf{a} | \mathbf{q}, \mathbf{o} = 0) = \frac{\pi_{\text{old}}(\mathbf{a} | \mathbf{q}) p(\mathbf{o} = 0 | \mathbf{q}, \mathbf{a})}{\sum_{\mathbf{A}} \pi_{\text{old}}(\mathbf{a} | \mathbf{q}) p(\mathbf{o} = 0 | \mathbf{q}, \mathbf{a})} \propto \pi_{\text{old}}(\mathbf{a} | \mathbf{q}) [1 - r(\mathbf{q}, \mathbf{a})], \\
r_{\mathbf{q}} \pi^{+}(\mathbf{a} | \mathbf{q}) + [1 - r_{\mathbf{q}}] \pi^{-}(\mathbf{a} | \mathbf{q}) &= \pi_{\text{old}}(\mathbf{a} | \mathbf{q}),
\end{aligned}$$

Finally we copy training objective Eq. 9 below.

$$\mathcal{L}_{(\mathbf{a}, \mathbf{q}, r) \sim \mathcal{D}}^{\text{NFT}}(\theta) = r \left[ -\log \frac{\pi_{\theta}^{+}(\mathbf{a} | \mathbf{q})}{\pi_{\text{old}}(\mathbf{a} | \mathbf{q})} \right] + (1 - r) \left[ -\log \frac{1 - r_{\mathbf{q}} \frac{\pi_{\theta}^{+}(\mathbf{a} | \mathbf{q})}{\pi_{\text{old}}(\mathbf{a} | \mathbf{q})}}{1 - r_{\mathbf{q}}} \right] \tag{16}$$

**Theorem B.1 (NFT for continuous rewards).** *Consider the training objective of Eq. 16 with arbitrary continuous rewards  $r(\mathbf{q}, \mathbf{a}) \in [0, 1]$ . Suppose we can accurately estimate  $r_{\mathbf{q}} = \mathbb{E}_{\pi_{\text{old}}(\mathbf{a} | \mathbf{q})} r(\mathbf{q}, \mathbf{a})$ . Assuming unlimited data and model capacity, the optimal solution satisfies*

$$\forall \mathbf{q}, \mathbf{a} : \quad \pi_{\theta^{*}}^{+}(\mathbf{a} | \mathbf{q}) = \pi^{+}(\mathbf{a} | \mathbf{q})$$

*Proof.*

$$\begin{aligned}
\nabla_{\theta} \mathcal{L}_{(\mathbf{a}, \mathbf{q}, r) \sim \mathcal{D}}^{\text{NFT}}(\theta) &= \nabla_{\theta} \left[ r(\mathbf{q}, \mathbf{a}) [-\log \pi_{\theta}^{+}(\mathbf{a}|\mathbf{q})] + (1 - r(\mathbf{q}, \mathbf{a})) \left[ -\log \frac{\pi_{\text{old}}(\mathbf{a}|\mathbf{q}) - r_{\mathbf{q}} \pi_{\theta}^{+}(\mathbf{a}|\mathbf{q})}{1 - r_{\mathbf{q}}} \right] \right] \\
&= \left[ r(\mathbf{q}, \mathbf{a}) \left[ -\frac{1}{\pi_{\theta}^{+}(\mathbf{a}|\mathbf{q})} \right] + (1 - r(\mathbf{q}, \mathbf{a})) \left[ \frac{r_{\mathbf{q}}}{\pi_{\text{old}}(\mathbf{a}|\mathbf{q}) - r_{\mathbf{q}} \pi_{\theta}^{+}(\mathbf{a}|\mathbf{q})} \right] \right] \nabla_{\theta} \pi_{\theta}(\mathbf{a}|\mathbf{q}) \\
&= \left[ \frac{r_{\mathbf{q}} \pi_{\theta}^{+}(\mathbf{a}|\mathbf{q}) - r(\mathbf{q}, \mathbf{a}) \pi_{\text{old}}(\mathbf{a}|\mathbf{q})}{\pi_{\theta}^{+}(\mathbf{a}|\mathbf{q}) [\pi_{\text{old}}(\mathbf{a}|\mathbf{q}) - r_{\mathbf{q}} \pi_{\theta}^{+}(\mathbf{a}|\mathbf{q})]} \right] \nabla_{\theta} \pi_{\theta}(\mathbf{a}|\mathbf{q})
\end{aligned}$$

When  $r_{\mathbf{q}} \pi_{\theta}^{+}(\mathbf{a}|\mathbf{q}) - r(\mathbf{q}, \mathbf{a}) \pi_{\text{old}}(\mathbf{a}|\mathbf{q}) = 0$ , that is

$$\pi_{\theta}^{+}(\mathbf{a}|\mathbf{q}) = \frac{r(\mathbf{q}, \mathbf{a}) \pi_{\text{old}}(\mathbf{a}|\mathbf{q})}{r_{\mathbf{q}}} = \pi^{+}(\mathbf{a}|\mathbf{q})$$

We have  $\nabla_{\theta} \mathcal{L}_{(\mathbf{a}, \mathbf{q}, r) \sim \mathcal{D}}^{\text{NFT}}(\theta) = 0$  constantly holds.  $\square$

## C EXPERIMENT DETAILS

**General training setup.** All algorithms are implemented based on the official DAPO codebase within the VeRL framework. We use a learning rate of 1e-6 with a linear warm-up schedule across all experiments. At each rollout step, we generate 16 answers for each of 512 sampled questions, then split the data into 16 mini-batches and train the policy network for 16 gradient steps. Models are trained for 320 rollout steps, totaling over 5,000 gradient steps. Unless otherwise specified, we follow DAPO’s default design choices, including dynamic data sampling, token-level loss normalization, and no KL regularization. For 7B training, we restrict context lengths to 4K and use 64 NVIDIA H100 GPUs. For 32B training, we restrict context lengths to 32K (DAPO) or 16K (others), and use 128-256 NVIDIA H100 GPUs.

**DAPO.** is a variant of GRPO that has achieved state-of-the-art AIME performance on 32B models. Our NFT implementation is adapted from the official DAPO codebase, based on the VeRL framework (Sheng et al., 2024). NFT inherits most of DAPO’s hyperparameters and design choices, including dynamic data sampling, token-level loss normalization, and no KL regularization. We adopt a faithful implementation of the official DAPO codebase, keeping all hyperparameters untouched.

**NFT.** Compared with DAPO, NFT modifies the context length to 16K for 32B training. We find this does not significantly affect performance, but noticeably reduces data collection time. Another difference is that we remove the overlong reward shaping technique of DAPO to ensure a binary reward outcome. In our setting, truncated answers are treated as negative, which sufficiently discourages overlong answers. The negative ratio clipping parameter is set to  $\epsilon = 1.0$ , and the prompt weight is defined as  $\omega(\mathbf{q}) = 1 - r_{\mathbf{q}}$ .

**RFT** is a simple but effective SL baseline that only fine-tunes LLMs on positive answers and throws away negative data. In our implementation, the main difference between RFT and NFT is that RFT zeros out negative data loss and keeps a constant prompt weight  $\omega(\mathbf{q}) = 1$  during training. For RFT, we zero out negative data loss in NFT implementation and keep a constant prompt weight  $\omega(\mathbf{q}) = 1$  during training.

**GRPO** does not adopt the dynamic sampling technique proposed by DAPO. Rather, it simply leverages all data for training, even though the gradient for all-positive or all-negative questions should be zeroed out automatically by the algorithm itself (Yu et al., 2025). Other DAPO-related techniques are kept, such as setting positive clipping parameter to a higher value  $\epsilon'_{+} = 0.28$ . Since GRPO requires less time for sampling data, we train GRPO models for 580+ rollout steps instead of 320 steps, which roughly takes the same training time compared with DAPO experiments.

**Dr. GRPO** is modified from our GRPO implementation. The only difference is the removal of the std normalization when computing group-normalized advantages.

**Iterative DPO.** Comparing DPO with other RL algorithms in a head-to-head fashion is difficult because DPO requires paired data to calculate the training objective, while we sample 16 unpaired

answers for each question. To solve this problem, we take the implementation of InfoNCA (Chen et al., 2024), a variation of the DPO algorithm that can handle  $K > 2$  responses per question:

$$\mathcal{L}_{(\mathbf{q}, \mathbf{a}^{1:K}, \mathbf{r}^{1:K}) \sim \mathcal{D}}^{\text{InfoNCA}}(\theta) = - \sum_{k=1}^K \left[ \frac{r^{(k)}}{\sum_{i=1}^K r^{(i)}} \log \frac{e^{\beta R_{\theta}(\mathbf{q}, \mathbf{a}^k)}}{\sum_{i=1}^K e^{\beta R_{\theta}(\mathbf{q}, \mathbf{a}^i)}} \right]$$

InfoNCA is guaranteed to become DPO algorithm in  $K = 2$  settings. We ablate  $\beta \in \{0.1, 0.01, 0.02\}$  and report the best averaged validation result. We find InfoNCA training to be unstable and add an SFT regularization term to the original loss function for stabilizing the training process.

**Validation details.** Validation is performed with a top-p value of 0.7. The temperature is set to 1.0 for 7B models and 0.6 for 32B models, with context lengths matching the training configuration. We use `math-verify` (Kydlicek, 2024) as the verifier during training validation, and `simpleRL` (Zeng et al., 2025) for final evaluation. The DAPO-17k benchmark consists solely of training questions whose ground truth answers are integers and includes both a prefix and a lastfix for each question. Accordingly, for validation on AIME and AMC questions, we adapt the prompt format to match the training pattern. For other benchmarks with non-integer answers, question prompts remain unmodified.

## D ADDITIONAL EXPERIMENT RESULTS

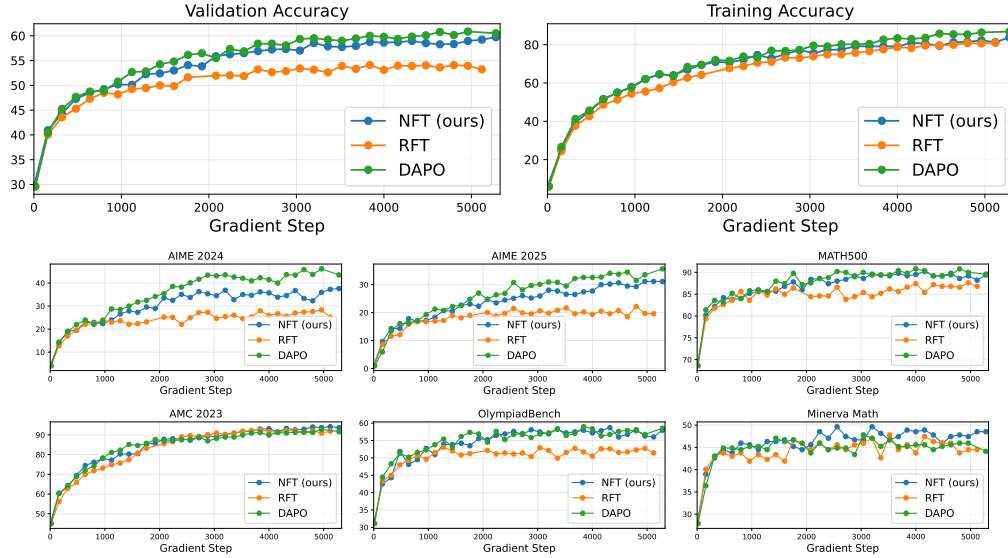


Figure 11: Training and validation accuracy curves for 32B experiments.