

# Fine-tuning LLMs with Synthetic Data

Pu Yang  
School of Mathematical Sciences, Peking University

2024.10.30

# Table of contents

- 1 Introduction
- 2 Technical 1: Multiple-step
  - RFT + Multiple-step
  - DPO + Multiple-step
- 3 Technical 2: Iterative Learning
- 4 Summary

# 1 Introduction

## 2 Technical 1: Multiple-step

- RFT + Multiple-step
- DPO + Multiple-step

## 3 Technical 2: Iterative Learning

## 4 Summary

# Introduction: Learning from Synthetic Data

- What is synthetic data?
  - ▶ Data generated by models rather than humans
- Why synthetic data?
  - ▶ SFT data is *expensive*, while synthetic data is much cheaper
  - ▶ SFT data is hard to make models *outperform* human (like imitation learning), while synthetic data allows models to explore
- How to learning form synthetic data? Two main question:
  - ▶ How to generate synthetic data? e.g. Teacher-student, **self-explore**
  - ▶ **How to use synthetic data for training?**

# Generate via Self-explore: Rejection Sampling Fine-Tuning (RFT)

- LLMs are tasked to self-generated *positive* data for rejection sampling fine-tuning (RFT).

Given a generative policy  $\pi$  and a binary reward function  $r(\mathbf{y}, \hat{\mathbf{y}}) \rightarrow \{0, 1\}$  which verifies if a new generation  $\hat{\mathbf{y}}$  is correct or not, A positive dataset  $\mathcal{D}_{\pi}^{+} = \{(\mathbf{x}, +\mathbf{y})\}$  where  $+\mathbf{y}$  is positive generated from  $\pi(\cdot | \mathbf{x})$ .

Then, we apply the next token prediction loss on  $\mathcal{D}_{\pi_{\text{sft}}}^{+}$

$$\max \mathbb{E}_{(\mathbf{x}, +\mathbf{y}) \sim \mathcal{D}_{\text{sft}}^{+}} [p(+\mathbf{y} | \mathbf{x})] \quad (1)$$

- ▶ e.g., code generation, where  $r$  is a Python interpreter.

# Generate via Self-explore: Direct Preference Optimization (DPO)

- LLMs are tasked to self-generated positive and negative data - which can together form a *pairwise* dataset for preference optimization, e.g. DPO.

A pairwise dataset  $\mathcal{D}_{\pi}^{\pm} = \{(x, +y, -y)\}$  where  $-y$  is negative generated from  $\pi(\cdot | x)$ .

Then we apply DPO loss on  $\mathcal{D}_{\pi}^{\pm}$

$$\min_{\pi} \mathcal{L}_{\text{DPO}}(\pi) := \mathbb{E}_{(x, +y, -y) \sim \mathcal{D}_{\text{sft}}^{\pm}} \left[ \sigma \left( \beta \log \frac{\pi(+y | x)}{\pi_{\text{sft}}(+y | x)} - \beta \log \frac{\pi(-y | x)}{\pi_{\text{sft}}(-y | x)} \right) \right] \quad (2)$$

- In general, performance: DPO > RFT > SFT > ICL

# Relation to Data Augmentation (DA)

- Traditional DA techniques aim at expanding the training dataset in a somewhat mechanical manner, e.g.
  - ▶ paraphrasing
  - ▶ back-translation
- Generating synthetic data with LLMs serves as an advanced DA method, which focuses on the generation of novel, context-rich training data tailored to specific domains and skills.

## 1 Introduction

## 2 Technical 1: Multiple-step

- RFT + Multiple-step
- DPO + Multiple-step

## 3 Technical 2: Iterative Learning

## 4 Summary



# RFT + Multiple-step<sup>1</sup>

Google DeepMind

2024-05-22

## Improve Mathematical Reasoning in Language Models by Automated Process Supervision

Liangchen Luo<sup>1\*</sup>, Yinxiao Liu<sup>1\*</sup>, Rosanne Liu<sup>1</sup>, Samrat Phatale<sup>1</sup>, Harsh Lara<sup>1</sup>, Yunxuan Li<sup>2</sup>, Lei Shu<sup>1</sup>, Yun Zhu<sup>1</sup>, Lei Meng<sup>2</sup>, Jiao Sun<sup>2</sup> and Abhinav Rastogi<sup>1</sup>

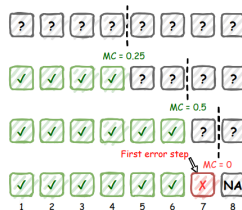
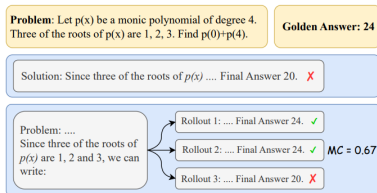
<sup>1</sup>Google DeepMind, <sup>2</sup>Google

<sup>1</sup>Liangchen Luo et al. "Improve Mathematical Reasoning in Language Models by Automated Process Supervision". In: [arXiv preprint arXiv:2406.06592](https://arxiv.org/abs/2406.06592) (2024).



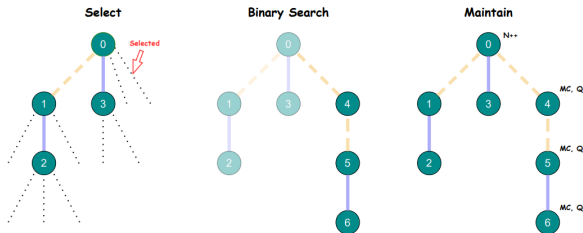
# Algorithm: Omega Process Reward Model (PRM)

- Build positive training examples via Monte Carlo Tree Search (MCTS)



(a) Monte Carlo estimation of a prefix solution.

(b) Error locating using binary search.



(c) Three stages in an iteration of the MCTS process.

# Algorithm: OmegaPRM

- Training with the positive training examples

$$\max \mathbb{E}_{(\mathbf{x}, \mathbf{Tree}) \sim \mathcal{D}_{\text{MCTS}}} [p(+\mathbf{y}_i \mid \mathbf{x}, +\mathbf{y}_{1:i-1})] \quad (3)$$

# DPO + Multiple-step<sup>2</sup>

---

## RL on Incorrect Synthetic Data Scales the Efficiency of LLM Math Reasoning by Eight-Fold

Amrith Setlur<sup>1</sup>, Saurabh Garg<sup>1</sup>, Xinyang (Young) Geng<sup>2</sup>, Naman Garg<sup>3</sup>, Virginia Smith<sup>1</sup> and Aviral Kumar<sup>2</sup>

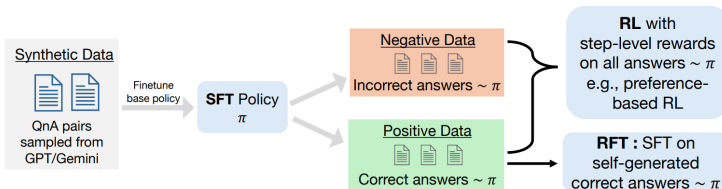
<sup>1</sup>Carnegie Mellon University, <sup>2</sup>Google DeepMind, <sup>3</sup>MultiOn

---

---

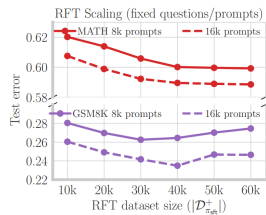
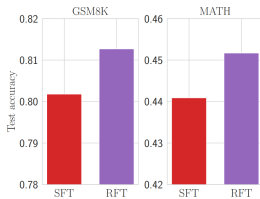
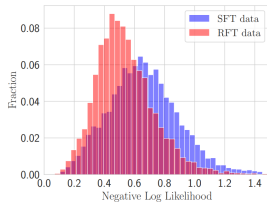
<sup>2</sup>Amrith Setlur et al. "RL on Incorrect Synthetic Data Scales the Efficiency of LLM Math Reasoning by Eight-Fold". In: [arXiv preprint arXiv:2406.14532](https://arxiv.org/abs/2406.14532) (2024).

# Overview



**Figure 1:** definitions for positive and negative synthetic data and how they are fed to SFT, RFT and step-level RL algorithms.

# Positive Data Improves Coverage, But Amplifies Spurious Correlations



- Under base LLM, RFT data sampled from  $\pi_{\text{sft}}$ , has higher likelihood than SFT data.
- RFT data with a single self-generated correct solution per problem outperforms SFT data of the same size.
- However, continuing to scale RFT data leads to test error saturation, or even worse test error (Contrary to the scaling law). Why? Incorrect/Irrelevant steps are not detected by our verifier.

# Negative Synthetic Data Enables Per-Step Credit Assignment

- Reasoning steps: The trace  $\mathbf{y}_i$  consists of several intermediate steps,  $\mathbf{y}_i = [\mathbf{y}_{i,1}, \dots, \mathbf{y}_{i,L}]$ .
- We formalize the notion of per-step credit using value functions from RL

$$Q_{\tilde{\pi}}(\underbrace{\mathbf{x}, \hat{\mathbf{y}}_{1:i-1}}_{\text{state}}; \underbrace{\hat{\mathbf{y}}_i}_{\text{action}}) = \underbrace{\mathbb{E}_{\mathbf{y}_{i+1:L}^{\text{new}} \sim \tilde{\pi}(\cdot | \mathbf{x}, \hat{\mathbf{y}}_{1:i})} [r([\hat{\mathbf{y}}_{1:i}, \mathbf{y}_{i+1:L}^{\text{new}}], \mathbf{y})]}_{\text{expected future reward under new actions (i.e., steps) sampled by policy } \tilde{\pi}} \quad (4)$$

- State: Problem  $\mathbf{x}$  and previous reasoning steps  $\hat{\mathbf{y}}_{1:i-1}$
  - Action: Current reasoning step  $\hat{\mathbf{y}}_i$
- Advantage of a given step  $\hat{\mathbf{y}}_i$ :

$$A_{\tilde{\pi}}(\mathbf{x}, \hat{\mathbf{y}}_{1:i-1}; \hat{\mathbf{y}}_i) = Q_{\tilde{\pi}}(\mathbf{x}, \hat{\mathbf{y}}_{1:i-1}; \hat{\mathbf{y}}_i) - Q_{\tilde{\pi}}(\mathbf{x}, \hat{\mathbf{y}}_{1:i-2}; \hat{\mathbf{y}}_{i-1}). \quad (5)$$

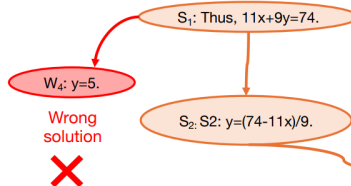
It is the gap between the Q-value of a state-action pair and the value function of the state.



# Illustration of Advantage Estimation from Negative Data on a Didactic Example

**Question:** 4 apples and 3 pears cost \$25, but 7 apples and 6 pears cost \$49. What is the cost of 1 apple?

Correct, but distracting step  
that may derail the model



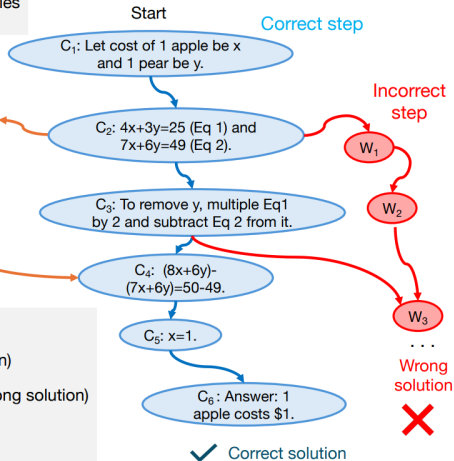
## Advantage function estimates

$A_{\pi}(C_2, C_2 \rightarrow S_1)$  is **low** (often leads to a wrong solution)

$A_{\pi}(C_2, C_2 \rightarrow C_3)$  is **midrange** (leads to correct and wrong solution)

$A_{\pi}(C_3, C_3 \rightarrow W_3)$  is **low** (leads to a wrong solution)

$A_{\pi}(C_3, C_3 \rightarrow C_4)$  is **high** (leads to the correct solution)



# Equivalence of Advantage-weighted RL and DPO with Per-Step Pairs

- Callback:

$$\min_{\pi} \mathcal{L}_{\text{DPO}}(\pi) := \mathbb{E}_{(\mathbf{x}, +\mathbf{y}, -\mathbf{y}) \sim \mathcal{D}_{\text{sft}}^{\pm}} \left[ \sigma \left( \beta \log \frac{\pi(+\mathbf{y} | \mathbf{x})}{\pi_{\text{sft}}(+\mathbf{y} | \mathbf{x})} - \beta \log \frac{\pi(-\mathbf{y} | \mathbf{x})}{\pi_{\text{sft}}(-\mathbf{y} | \mathbf{x})} \right) \right]$$

- We have the following theorem:

**Theorem 6.1** (Equivalence of advantage-weighted RL and DPO with per-step pairs). *The optimal policy from Equation 1 with  $\mathcal{D}_{\pi_{\text{sft}}}^{\pm}$  given by  $(\mathbf{x}, [\mathbf{y}_{1:i}, +\mathbf{y}_{i+1}], [\mathbf{y}_{1:i}, -\mathbf{y}_{i+1}])$  where the positive and negative traces share prefix  $\mathbf{y}_{1:i} \sim \pi_{\text{sft}}$ , and  $-\mathbf{y}_{i+1} \sim \pi_{\text{sft}}(\cdot | \mathbf{x}, \mathbf{y}_{1:i})$ ,  $+\mathbf{y}_{i+1} \sim \sigma(A_{\tilde{\pi}}(\mathbf{x}, \mathbf{y}_{1:i}; \cdot) - A_{\tilde{\pi}}(\mathbf{x}, \mathbf{y}_{1:i}; -\mathbf{y}_{i+1}))$ , is identical to the optima of the advantage-weighted RL objective:*

$$\max_{\pi} \mathbb{E}_{\mathbf{x} \sim p_{\text{syn}}(\mathbf{x}), \mathbf{y} \sim \pi_{\text{sft}}(\cdot | \mathbf{x})} \left[ \sum_{i=1}^L \log \pi(\mathbf{y}_i | \mathbf{x}, \mathbf{y}_{0:i-1}) \cdot \exp(A_{\tilde{\pi}}(\mathbf{x}, \mathbf{y}_{0:i-1}; \mathbf{y}_i) / \beta) \right]. \quad (4)$$

# Per-step DPO Algorithm

---

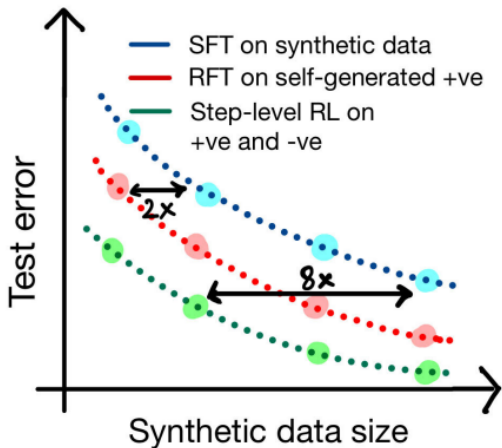
**Algorithm 1** Per-step DPO (Part 1: Practical version for most experiments; Parts 1 + 2: Complete version)

---

**Require:** Synthetic dataset:  $\mathcal{D}_{\text{syn}}$ , SFT policy finetuned on  $\mathcal{D}_{\text{syn}}$ :  $\pi_{\text{sft}}$ , sampling policy  $\tilde{\pi}$ .

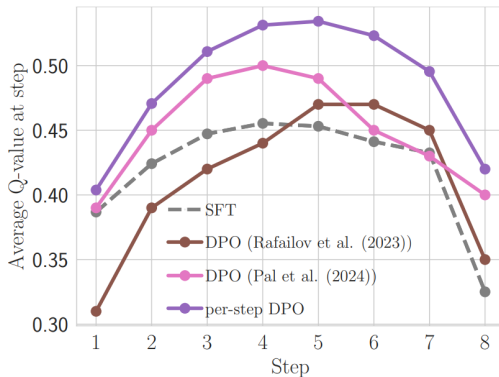
- 1: Initialize per-step DPO dataset  $\mathcal{D}_{\pi_{\text{sft}}}^{\pm} \leftarrow \{\}$ .
  - 2: **for**  $(x, y) \in \mathcal{D}_{\text{syn}} \cup \mathcal{D}_{\pi_{\text{sft}}}^{\pm}$  **do**
  - 3:   # **Part 1: Identify critical steps in incorrect responses**
  - 4:   Sample multiple incorrect answers  $-\hat{y} \sim \pi_{\text{sft}}(\cdot | x)$ , and collect them in set  $\mathcal{C}(x)$ .
  - 5:   **for**  $-\hat{y} := [-\hat{y}_1, \dots, -\hat{y}_L] \in \mathcal{C}(x)$  **do**
  - 6:     Compute the Monte Carlo estimate for  $Q_{\tilde{\pi}}(x, -\hat{y}_{1:i-1}; -\hat{y}_i)$  for each step  $-\hat{y}_i$ .
  - 7:     If  $-\hat{y}_c$  is the first step with least  $Q_{\tilde{\pi}}(x, -\hat{y}_{1:i-1}; -\hat{y}_i)$ , then  $\mathcal{D}_{\pi_{\text{sft}}}^{\pm} \leftarrow \mathcal{D}_{\pi_{\text{sft}}}^{\pm} \cup \{(x, y, -\hat{y}_{1:c})\}$ .
  - 8:   **end for**
  - 9:   # **Part 2: Identify spurious steps in correct responses**
  - 10:   Sample multiple correct answers  $+\hat{y} \sim \pi_{\text{sft}}(\cdot | x)$ , and collect them in set  $\mathcal{C}'(x)$ .
  - 11:   **for**  $+\hat{y} := [+\hat{y}_1, \dots, +\hat{y}_L] \in \mathcal{C}'(x)$  **do**
  - 12:     Compute the Monte Carlo estimate for  $Q_{\tilde{\pi}}(x, +\hat{y}_{1:i-1}; +\hat{y}_i)$  for each step  $+\hat{y}_i$ .
  - 13:     If  $+\hat{y}_c$  is the first step with least  $Q_{\tilde{\pi}}(x, +\hat{y}_{1:i-1}; +\hat{y}_i)$ , then  $\mathcal{D}_{\pi_{\text{sft}}}^{\pm} \leftarrow \mathcal{D}_{\pi_{\text{sft}}}^{\pm} \cup \{(x, y, +\hat{y}_{1:c})\}$ .
  - 14:   **end for**
  - 15: **end for**
  - 16: Optimize DPO loss in Equation (1) on  $\mathcal{D}_{\pi_{\text{sft}}}^{\pm}$  with  $\pi_{\text{sft}}$  as the reference policy.
-

# Results



- Per-step DPO > RFT > SFT

# Ablation: Number of Steps



- When the number of steps is small, multi-step reasoning is treated as one step.
- When the number of steps is large, the Monte Carlo estimation may not be accurate.

# Example: Advantage Estimation

## Advantage Estimates Computed Over a Positive Model Generated Sample

### Question:

Betty wants a new wallet which costs \$100. Betty has only half of the money she needs. Her parents give her \$15 for that purpose, and her grandparents twice as much as her parents. How much more money does Betty need?

### Positive response with identified spurious step:

Betty's parents gave her 15, so her grandparents gave her  $2 * 15 = 15 * 2 = 30$ . In total, Betty received  $30 + 15 = 30 + 15 = 45$ . The total amount of money Betty needs is  $100 * 2 = 100 * 2 = 200$ . Betty needs  $100 / 2 = 100 / 2 = 50$  in total. Betty still needs to save  $50 - 45 = 50 - 45 = 5$ . The answer is 5

## 1 Introduction

## 2 Technical 1: Multiple-step

- RFT + Multiple-step
- DPO + Multiple-step

## 3 Technical 2: Iterative Learning

## 4 Summary

---

# Monte Carlo Tree Search Boosts Reasoning via Iterative Preference Learning

---

**Yuxi Xie<sup>1\*</sup> Anirudh Goyal Wenyue Zheng<sup>1</sup> Min-Yen Kan<sup>1</sup>**  
**Timothy Lillicrap<sup>2</sup> Kenji Kawaguchi<sup>1</sup> Michael Shieh<sup>1</sup>**

<sup>1</sup> National University of Singapore

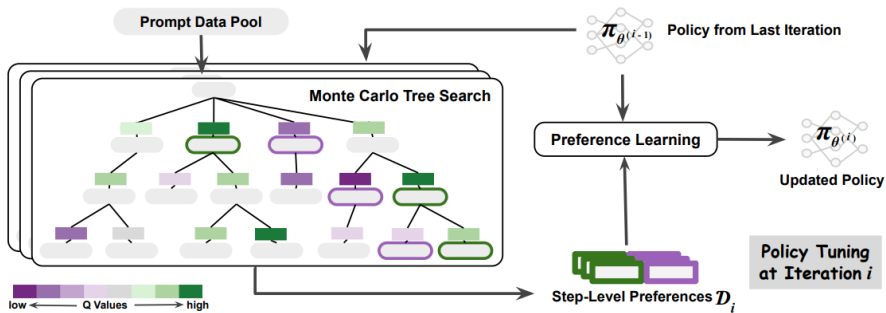
<sup>2</sup> Google DeepMind

---

<sup>3</sup>Yuxi Xie et al. "Monte Carlo Tree Search Boosts Reasoning via Iterative Preference Learning". In: [arXiv preprint arXiv:2405.00451](https://arxiv.org/abs/2405.00451) (2024).

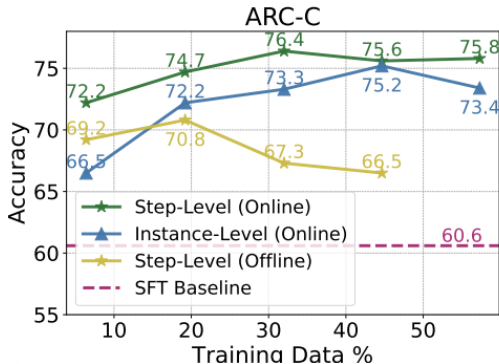


# Overview



- Use action values  $Q$  estimated by MCTS to assign the preferences

## Main Results



- Offline setting can fail with high probability if the sampling policy differs too much from the current policy.
- we can indeed avoid this failure case in the online setting.

- 1 Introduction
- 2 Technical 1: Multiple-step
  - RFT + Multiple-step
  - DPO + Multiple-step
- 3 Technical 2: Iterative Learning
- 4 Summary

## Summary: Exploitation and Exploration

We talk about how to learn from synthetic data.

- Explore more efficiently
- Make full use of synthetic data