

# GPO: Learning from Critical Steps to Improve LLM Reasoning

Author: Jiahao Yu, Zelei Cheng, Xian Wu, Xinyu Xing

Presenter: Jiahao Yu

Northwestern

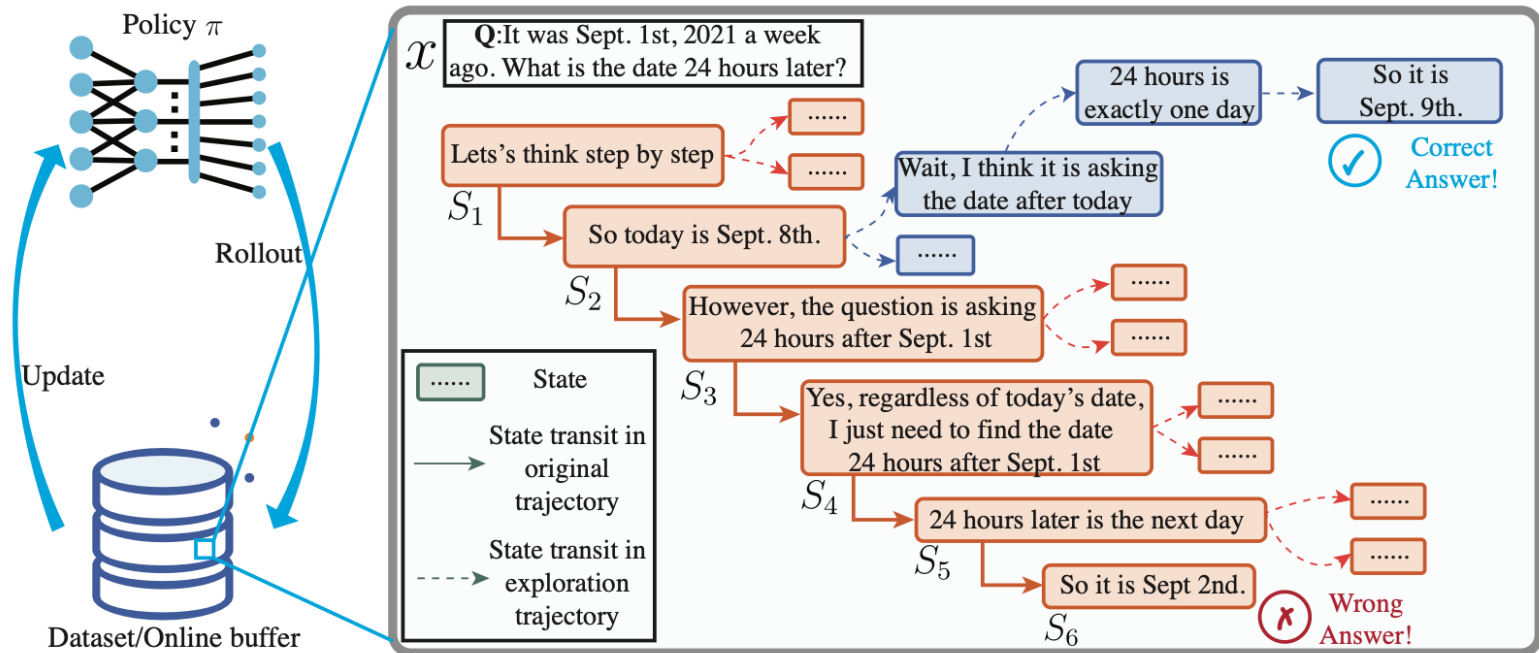
 Meta

 Capital One

# Motivation

1. LLM Reasoning can be brittle
2. Standard fine-tuning treat reasoning trajectories as a whole
3. Previous method fail to pinpoint which specific step was critical

# Our Solution: GPO



# A General Framework

---

**Algorithm 1:** GPO Optimization Framework

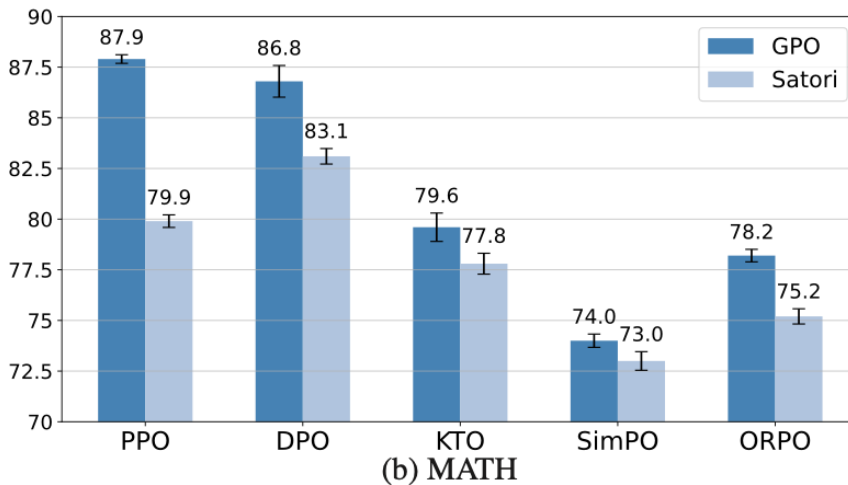
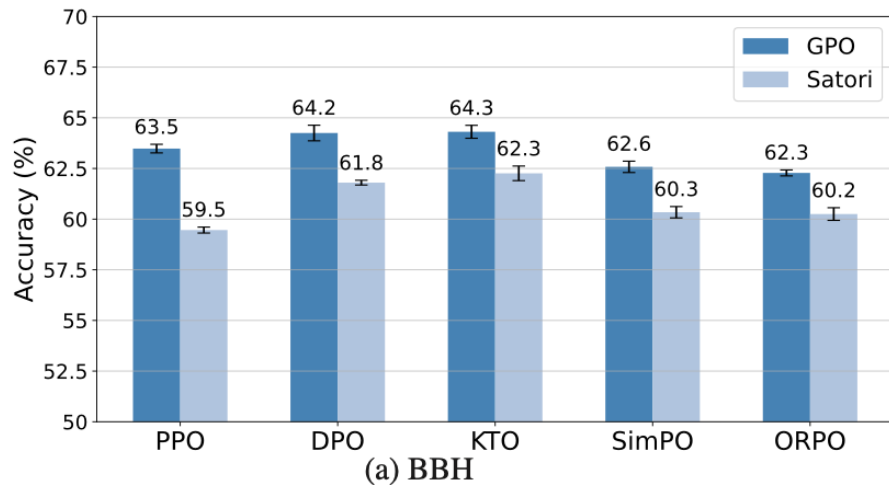
---

- 1: **Procedure-I: Online Policy Training (PPO-based)**
  - 2: **Input:** Initial LLM policy  $\pi^0 = \pi_{\text{ref}}$ , reasoning dataset  $D_r$
  - 3: **for** iteration = 1, 2, ...,  $T$  **do**
  - 4:    $\mathcal{D} \leftarrow \emptyset$
  - 5:   **for** n = 1, 2, ... **do**
  - 6:     Random sampling a question  $x$  from the reasoning dataset  $D_r$
  - 7:     Run  $\pi^t$  to generate a  $K$ -step reasoning trajectory  $y = (y_0, \dots, y_{K-1})$  split by newlines
  - 8:     Identify the critical step  $y_m$  with maximal advantage  $A^{\pi^t}(x, y_{0:i-1}; y_i)$
  - 9:     Reset  $\pi^t$  to  $y_m$  and roll-out  $\pi^t$  to generate trajectory  $y' = (y_m, \dots, y'_{K-1})$
  - 10:    Add trajectory  $y'$  and the final reward  $r$  to  $\mathcal{D}$
  - 11:   **end for**
  - 12:   Optimize  $\pi^t$  with respect to the policy gradient loss (e.g., PPO loss) in **Eqn. 1** on  $\mathcal{D}$
  - 13: **end for**
  - 14: **Procedure-II: Preference Data Generation and Optimization (DPO-based)**
  - 15: **Input:** Supervised-finetuned base policy  $\pi_{\text{ref}}$ , reasoning dataset  $D_r$
  - 16:  $\mathcal{D} \leftarrow \emptyset$
  - 17: **for** iteration = 1, 2, ...,  $T$  **do**
  - 18:   Repeat the sampling, trajectory generation using  $\pi_{\text{ref}}$ , and critical step identification as in **Procedure-I** to extract the important step  $y_m$  from trajectory  $y$ .
  - 19:   Generate two continuations starting from  $y_m$  to obtain a positive trajectory  $y^+ = (y_0, \dots, y_m, \dots, y_{K-1}^+)$  and a negative trajectory  $y^- = (y_0, \dots, y_m, \dots, y_{K-1}^-)$
  - 20:   Add the preference pair  $(x, y^+, y^-)$  to  $\mathcal{D}$
  - 21: **end for**
  - 22: Optimize  $\pi$  with respect to the preference loss (e.g., DPO loss) in **Eqn. 2** on  $\mathcal{D}$
-

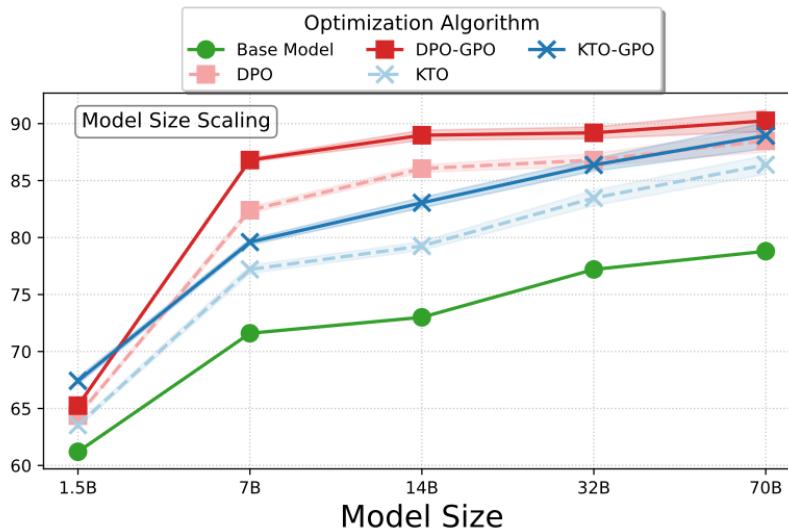
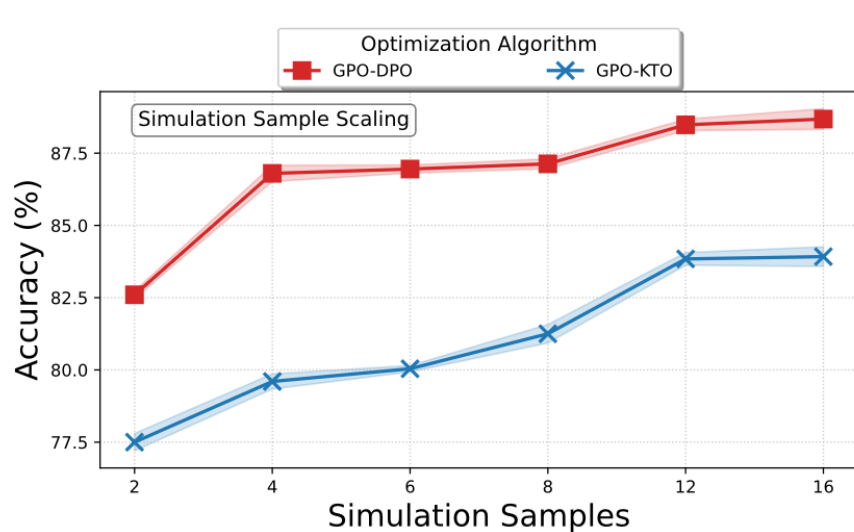
# Experimental Results

Algorithms	Test Accuracy (%)						
	BBH	MATH	GSM8K	MMLU	MMLUPro	AIME-2024	AIME-2025
Base Model	59.97	71.60	86.50	54.09	38.80	13.33	16.67
PPO	61.82	79.60	86.96	56.66	47.47	26.67	23.33
GPO-PPO	<b>63.48</b>	<b>87.80</b>	<b>87.44</b>	<b>59.39</b>	<b>51.05</b>	<b>30.00</b>	<b>26.67</b>
DPO	63.20	82.40	86.05	57.08	48.28	20.00	20.00
GPO-DPO	<b>64.25</b>	<b>86.80</b>	<b>88.48</b>	<b>58.93</b>	<b>51.93</b>	<b>26.67</b>	<b>26.67</b>
KTO	62.86	77.20	89.31	59.42	49.02	20.00	20.00
GPO-KTO	<b>64.31</b>	<b>79.60</b>	<b>90.25</b>	<b>61.35</b>	<b>50.52</b>	<b>23.33</b>	<b>26.67</b>
SimPO	61.97	72.20	86.58	56.93	45.70	20.00	23.33
GPO-SimPO	<b>62.58</b>	<b>74.00</b>	<b>88.35</b>	<b>57.44</b>	<b>47.74</b>	<b>23.33</b>	<b>26.67</b>
ORPO	61.75	75.20	87.26	57.72	46.66	20.00	20.00
GPO-ORPO	<b>62.28</b>	<b>78.20</b>	<b>88.17</b>	<b>58.72</b>	<b>48.65</b>	<b>23.33</b>	<b>23.33</b>

# Ablation Study



# Scaling Behavior



# Human Alignment

A user study found a strong correlation. Steps GPO identified as 'critical' were also chosen by human evaluators **44-88%** of the time as the most pivotal moment in the reasoning process.

# Thank you!

Code: <https://github.com/sherdencooper/GPO>