# **PGPO:** Enhancing Agent Reasoning via Pseudocode-style Planning Guided Preference Optimization

Zouying Cao<sup>1,3,4</sup>, Runze Wang<sup>2</sup>, Yifei Yang<sup>1,3,4</sup>, Xinbei Ma<sup>1,3,4</sup>, Xiaoyong Zhu<sup>2</sup>, Bo Zheng<sup>2,\*</sup>, Hai Zhao<sup>1,3,4,\*</sup>

<sup>1</sup>School of Computer Science, Shanghai Jiao Tong University, <sup>2</sup>Taobao & Tmail Group of Alibaba

<sup>3</sup>Key Laboratory of Shanghai Education Commission for Intelligent Interaction and Cognitive Engineering, Shanghai Jiao Tong University, <sup>4</sup>Shanghai Key Laboratory of Trusted Data Circulation and Governance in Web3

### Introduction

- Due to the dominance of natural language (NL) in agent reasoning, existing researches mainly focus on generating NL plans. We investigate the effectiveness of **pseudocode-style plans (P-code Plan)** in agent reasoning, which are more concise and structured than NL plans. P-code Plan demonstrates its superiority in boosting the generalization ability of LLM agents.
- We further propose a pseudocode-style <u>Planning Guided Preference Optimization</u> method named **PGPO** which empowers LLM agents with enhanced reasoning capabilities under the guidance of pseudocode-style plans. Analyses reveal the advantage of PGPO in reducing action errors and omissions during reasoning.





# Methodology: P-code Plan-Guided Agent Learning

#### • Contact us: Expert Dataset **Contrastive Trajectory Datasets** D. Explo zouyingcao@sjtu.edu.cn with P-code Plan Plan-driven Score yunze.wrz@alibaba-inc.com Reward Expert zhaohai@cs.sjtu.edu.cn क्रिके Trajector $\triangleleft$ DPO Loss $\mathcal{L}_p$ + DPO Loss $\mathcal{L}_f$ Plan-following 🄇 Full paper + SFT Loss Ls Reward O Observat Preference Thought Optimization base agent $\pi_{\theta}$ (a) Supervised (c) Planning-guided Agent (b) Exploration Stage for Planning-oriented Trajectory Collection **Fine Tuning Stage Learning Stage**

### **Experiments**

(1) Main	n Res	ults.					new sta	te-of-the-a	urt perforn	nance
Method	Llama-2-7B					Llama-2-13B				
	ALI	World UnSeen	WebShop	TextCraft	Avg.	ALI	F <b>World</b> UnSeen	WebShop	TextCraft	Avg.
SFT	60.0	67.2	60.2	28.0	53.9	67.1	67.9	62.2	29.0	56.6
ETO	68.6	72.4	67.4	35.0	60.9	75.0	69.4	68.9	42.0	63.8
IPR	<u>70.3</u>	74.7	71.3	34.0	<u>62.6</u>	75.0	76.9	72.2	39.0	<u>65.8</u>
PGPO	76.4	76.9	72.2	43.0	67.1	77.1	77.6	73.7	48.0	69.1
			Llama-3-8	В		Mistral-7B				
Method	ALI Seen	F <b>World</b> UnSeen	WebShop	TextCraft	Avg.	ALI Seen	F <b>World</b> UnSeen	WebShop	TextCraft	Avg.
SFT	67.1	72.4	61.2	20.0	55.2	72.1	68.7	61.8	31.0	58.4
ETO	72.1	73.1	66.2	36.0	61.9	75.0	72.4	66.2	<u>38.0</u>	62.9
IPR	72.9	<u>73.9</u>	72.0	<u>38.0</u>	<u>64.2</u>	<u>73.6</u>	<u>73.1</u>	69.6	36.0	<u>63.1</u>
PGPO	75.0	76.9	72.3	46.0	67.6	75.0	77.6	<u>69.0</u>	45.0	66.7

Table 1: Main results of PGPO compared to training-based baselines on ALFWorld, WebShop and TextCraft.

#### alleviate the need for few-shot context achieve comparable or even better results

Method	ALI Seen	F <b>World</b> UnSeen	WebShop	TextCraft	
ReAct+GPT-4	42.9	38.1	63.2	28.0	
ReAct+GPT-3.5	7.9	10.5	62.4	20.0	
ADaPT+GPT-4	75.0	69.4	64.8	48.0	
ADaPT+GPT-3.5	70.3	<u>71.6</u>	62.7	26.0	
PGPO+Llama-2-7B	76.4	76.9	72.2	43.0	
PGPO+Llama-3-8B	<u>75.0</u>	76.9	72.3	<u>46.0</u>	

Table 2: Comparative experiments vs. prompt-based baselines

#### **(2)** Approach Ablations.

	ALF	World	Walchan	T	
	Seen	UnSeen	webShop	lexicran	
PGPO	76.4	76.9	72.2	43.0	
- P-code	75.74 0.7	71.64 5.3	69.6↓ <u>2.6</u>	40.01 3.0	
- L <sub>f</sub>	74.3 1 2.1	75.4 1.5	70.4 1.8	41.04 2.0	
- L.	69.31 7.1	68.71 8.2	64.817.4	35.01 8.0	

Table 3: Approach ablations of PGPO. - *P-code* represents using NL plans to replace P-code Plans. -  $L_f$  denotes leaving out the estimation of plan-following reward, followed by the removal of  $L_f$ . -  $L_g$  means the removal of SFT loss.

# **(3)** Ablation on Optimization Iterations.



Figure 3: Ablation study on optimization iterations. (a) provides a comparison of the performance of PGPO against ETO and IPR across varying iterations. (b) shows the influence of increasing iterations on PGPO across different base models. iter=0 is the SFT stage.

## P-code Plan Generation Pipeline

- **Thought Extraction**: extract agent thoughts from existing ReAct-style expert trajectories, which involves task-specific knowledge for planning.
- **Plan Distillation**: instruct powerful LLM to summarize the step-by-step plan following the predefined P-code Plan format via few-shot prompting strategy.
- Plan Verification: minor manual refinement is needed to guarantee the quality of the LLM-generated P-code Plans.



Figure 2: Overview of P-code Plan generation pipeline. We first extract the thought part from existing ReActstyle datasets. Then, we prompt GPT-40 to summarize the thought process into high-level plans. Pseudocodestyle plans are finally structured with predefined formats, followed by manual verification to ensure accuracy.



ŝ

# (4) Analysis on training time efficiency.

PGPO delivers a 9% performance improvement while maintains reasonable training efficiency, requiring less than twice the time cost of ETO.

# (5) P-code Plan guidance can reduce the incidence of action errors and omissions in reasoning.



<sup>6</sup> Step-wise reward does not necessarily elicit better LLM agents.



It is still challenging to accurately determine the contribution of the intermediate step, thus introducing step-wise reward instead plays a negative role in agent reasoning.

### Conclusion

- P-code Plan can capture efficient structural logic of reasoning compared with NL plans, suitable for LLM agent's generalization to analogous agent tasks.
- With two planning-oriented rewards, PGPO further enhances LLM agents' ability to generate high-quality P-code Plans and subsequent reasoning.

