



Iterative Tool Usage Exploration for Multimodal Agents via Step-wise Preference Tuning

Pengxiang Li^{1,2*}, Zhi Gao^{2,3*}, Bofei Zhang², Yapeng Mi², Xiaojian Ma²,
Chenrui Shi^{1,2}, Tao Yuan², Yuwei Wu^{1,4} ✉, Yunde Jia^{4,1}, Song-Chun
Zhu^{2,3,5}, Qing Li² ✉

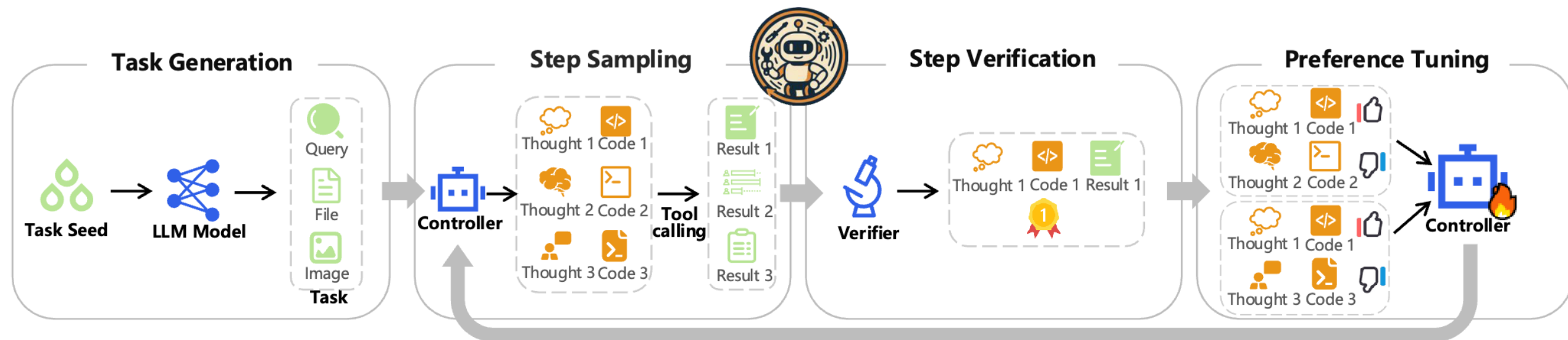
¹ Beijing Institute of Technology ² Beijing Institute of General Artificial Intelligence (BIGAI) ³ Peking University ⁴ Shenzhen MSU-BIT University
⁵ Tsinghua University



WeChat

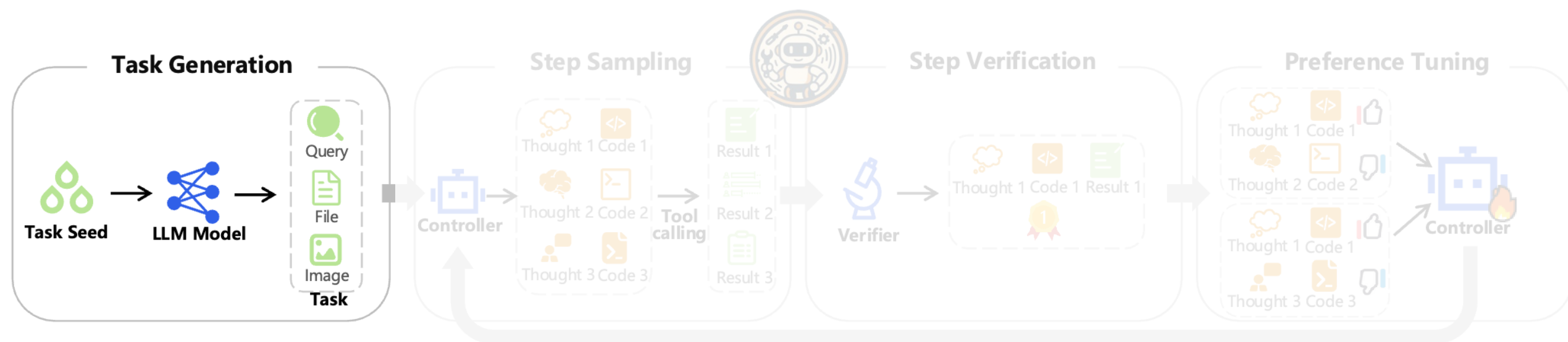


Tool Usage Exploration



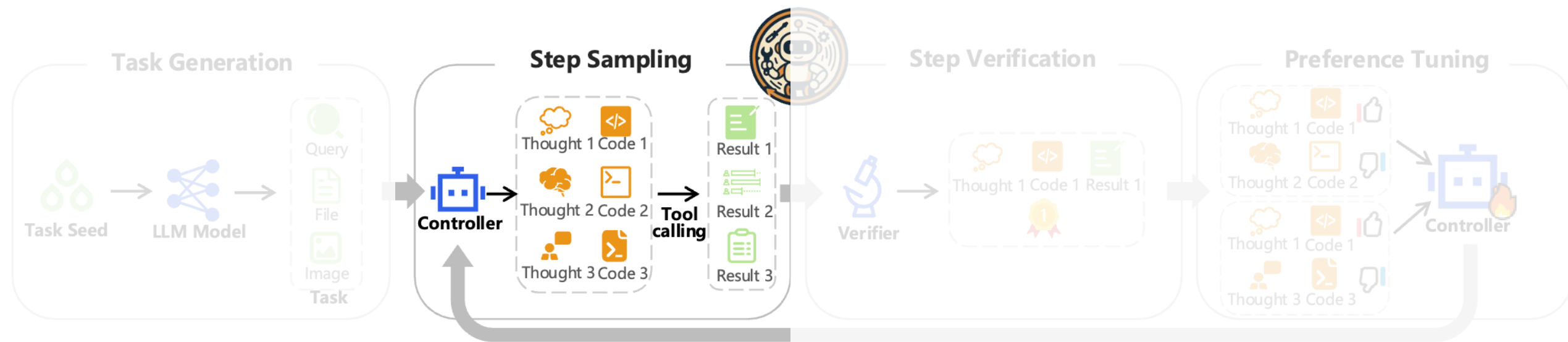
This paper introduces an online self-exploration loop that enables multimodal agents to self-improve via AI-generated tasks and LLM-verified preference tuning **without human annotations**.

Tool Usage Exploration



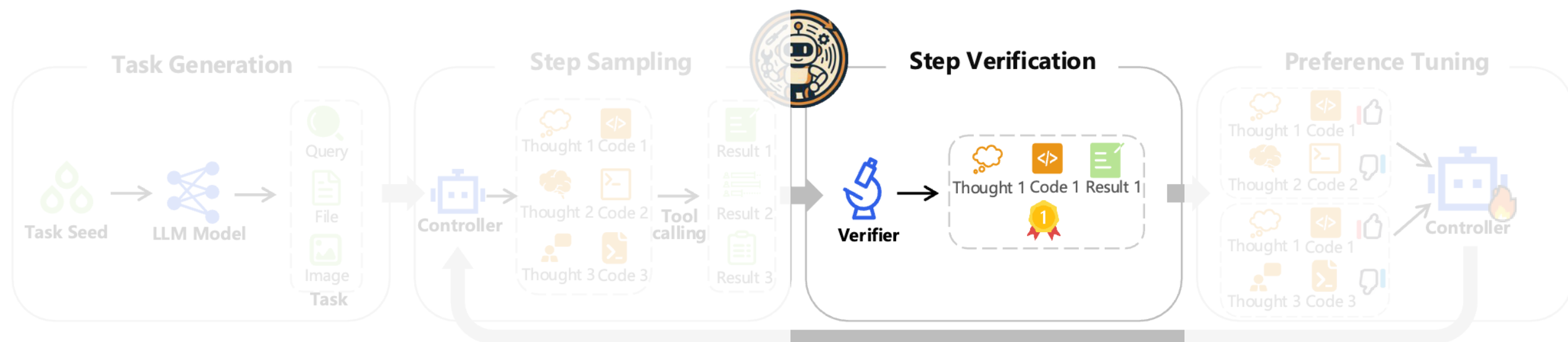
Task generalization: Use some seed to generate tasks.

Tool Usage Exploration



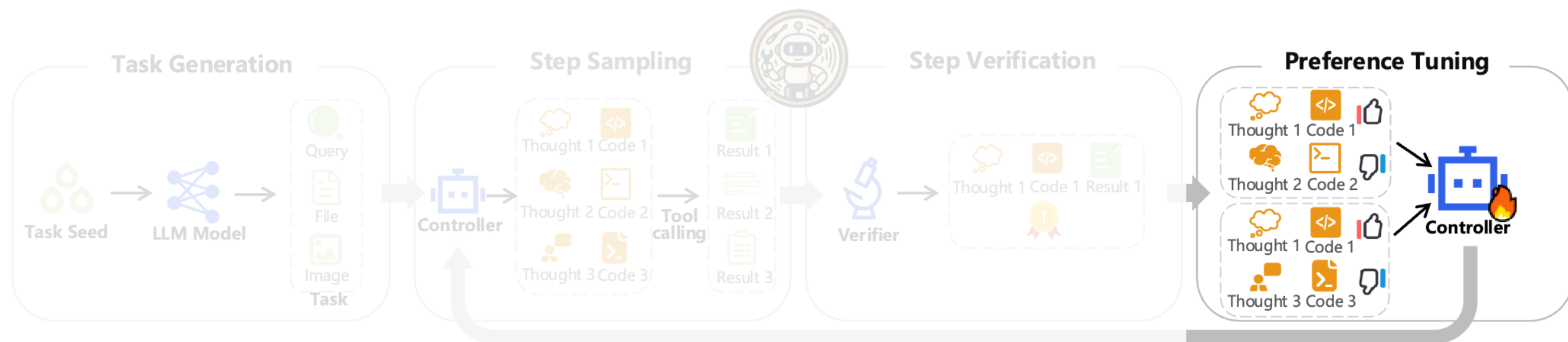
Step Sampling: Sampling possible solutions in one step.

Tool Usage Exploration



Step Verification: Use AI feedback to rank the sampled solutions.

Tool Usage Exploration



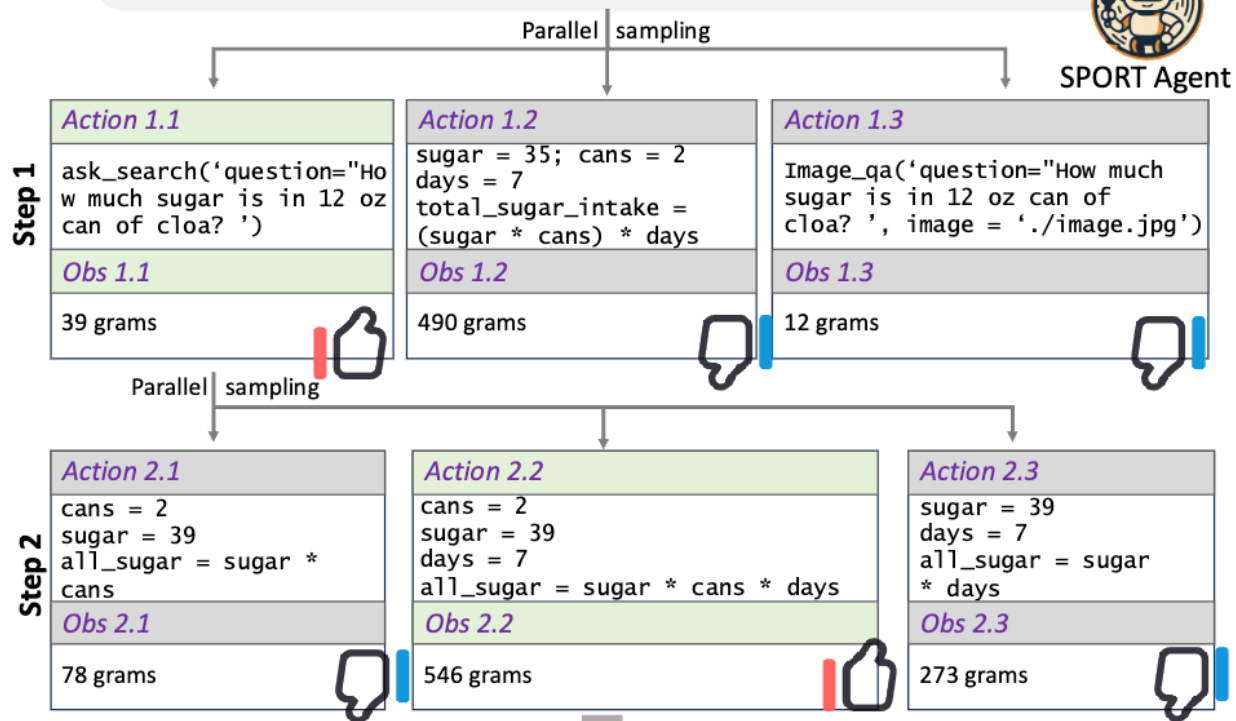
Preference Tuning: Perform step-wise DPO

$$\mathcal{L}(\theta) = -\mathbb{E}_{(x_i, a_i^{pre}, a_i^{dis}) \sim \mathcal{D}} [\log \sigma(\beta \log \frac{\pi_{\theta}(a_i^{pre} | x_i)}{\pi_{ref}(a_i^{pre} | x_i)} - \beta \log \frac{\pi_{\theta}(a_i^{dis} | x_i)}{\pi_{ref}(a_i^{dis} | x_i)})],$$

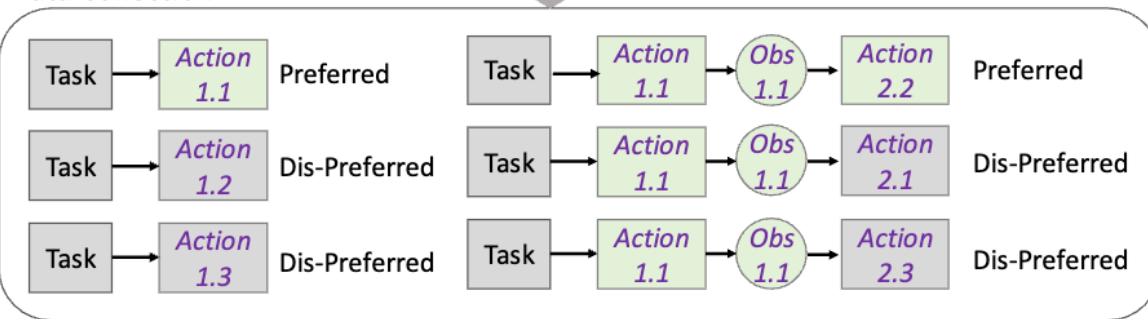
Step-wise Preference Data



Task: How many grams of sugar will I take in if I drink two drinks of this 12 oz can a day like that in the picture for a week?



Data Collection



Results on GTA benchmark

Method	Controller	GTA			GAIA			
		ToolAcc	CodeExec	AnsAcc	Level 1	Level 2	Level 3	AnsAcc
Closed-source Controller								
Lego Agent	GPT-4	-	-	46.59	-	-	-	-
Lego Agent	GPT-4o	-	-	41.52	-	-	-	-
Warm-up Agent	GPT-4-turbo	-	-	-	30.20	15.10	0.00	17.60
HF Agent	GPT-4o	63.41	95.12	57.05	47.17	31.40	11.54	33.40
HF Agent	GPT-4o	63.41	95.12	57.05	47.17	31.40	11.54	33.40
Compared with SFT agents. Higher AnsAcc (6.41%↑) on GTA, Higher AnsAcc (3.64%↑) on GAIA.								
HF Agent	LLaVA-Next-8B	14.97	23.08	14.10	9.43	1.16	0.00	3.64
HF Agent	InternVL2-8B	36.75	52.18	32.05	7.55	4.65	0.00	4.85
HF Agent	MiniCPM-V-8.5B	36.59	56.10	33.97	13.21	5.81	0.00	7.27
HF Agent	Qwen2-VL-7B	44.85	65.19	42.31	16.98	8.14	0.00	9.70
T3-Agent	MAT-MiniCPM-V-8.5B	65.85	80.49	52.56	26.42	11.63	3.84	15.15
T3-Agent	MAT-Qwen2-VL-7B	64.63	84.32	53.85	26.42	15.12	3.84	16.97
Ours								
SPORT Agent	Tuned-Qwen2-VL-7B	72.41	91.87	60.26	35.85	16.28	3.84	20.61

Results on GTA and GAIA benchmark

Method	Controller	GTA			GAIA			
		ToolAcc	CodeExec	AnsAcc	Level 1	Level 2	Level 3	AnsAcc
Closed-source Controller								
Lego Agent	GPT-4	-	-	46.59	-	-	-	-
Lego Agent	GPT-4o	-	-	41.52	-	-	-	-
Warm-up Agent	GPT-4-turbo	-	-	-	30.20	15.10	0.00	17.60
HF Agent	GPT-4o	63.41	95.12	57.05	47.17	31.40	11.54	33.40
HF Agent	GPT-4o-mini	56.10	100.00	57.69	33.96	27.91	3.84	26.06
Open-Source Controller								
HF Agent	LLaVA-NeXT-8B	14.97	25.08	14.10	9.43	1.16	0.00	3.64
HF Agent	InternVL2-8B	36.75	52.18	32.05	7.55	4.65	0.00	4.85
HF Agent	MiniCPM-V-8.5B	36.59	56.10	33.97	13.21	5.81	0.00	7.27
HF Agent	Qwen2-VL-7B	44.85	65.19	42.31	16.98	8.14	0.00	9.70
T3-Agent	MAT-MiniCPM-V-8.5B	65.85	80.49	52.56	26.42	11.63	3.84	15.15
T3-Agent	MAT-Qwen2-VL-7B	64.63	84.32	53.85	26.42	15.12	3.84	16.97
Ours								
SPORT Agent	Tuned-Qwen2-VL-7B	72.41	91.87	60.26	35.85	16.28	3.84	20.61

Examples

