



Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin, Liwen Zhang, Zhenglin Wang, Zhengwei Tao, Dingchu Zhang, Zekun Xi, Xiangru Tang, Yong Jiang, Pengjun Xie, Fei Huang, Jingren Zhou

DeepResearch Team, Tongyi Lab, Alibaba Group

wujialongml@gmail.com



Introduction

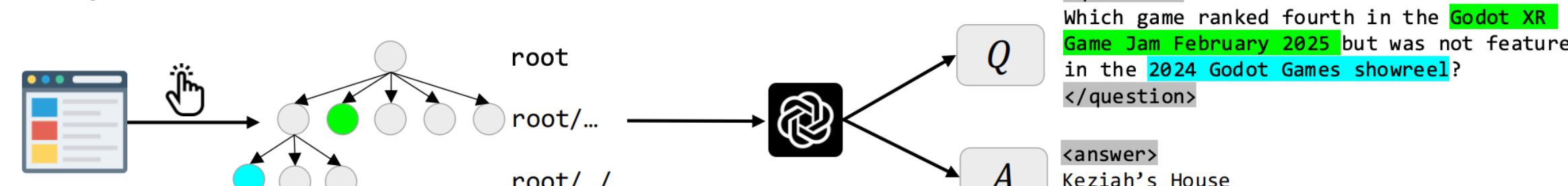
Deep Research has demonstrated strong deep information seeking capabilities through end-to-end reinforcement learning (RL) training.

How to build a web agent like Deep Research from scratch?

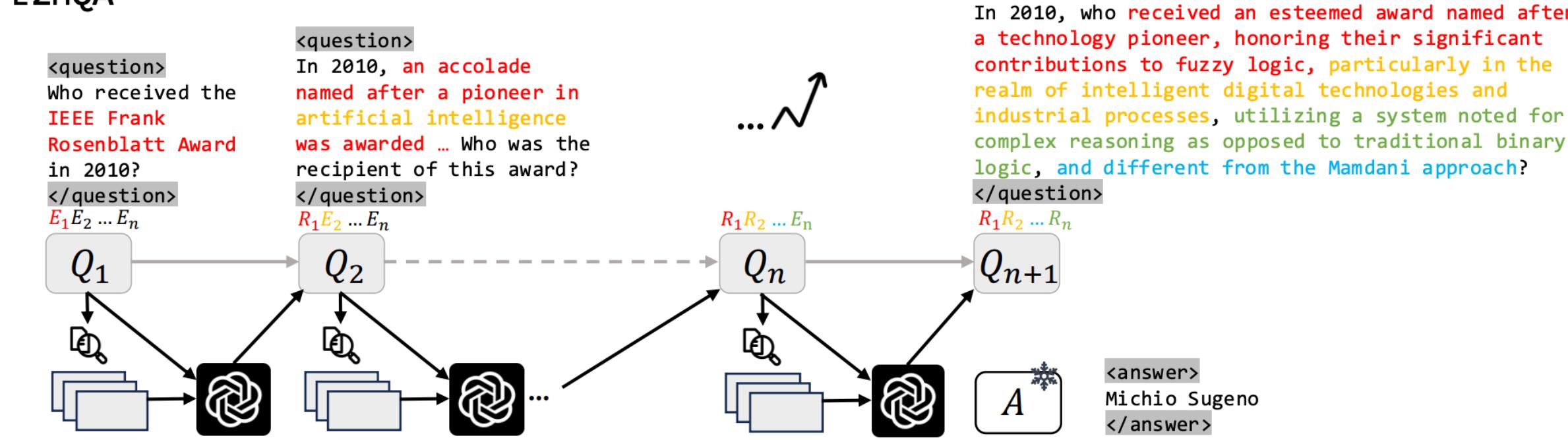
- (1) Acquiring high-quality synthetic QA pairs
- (2) Constructing reliable trajectories that support long-horizon reasoning and task decomposition
- (3) Agentic training strategies

Data

CRAWLQA



E2HQA

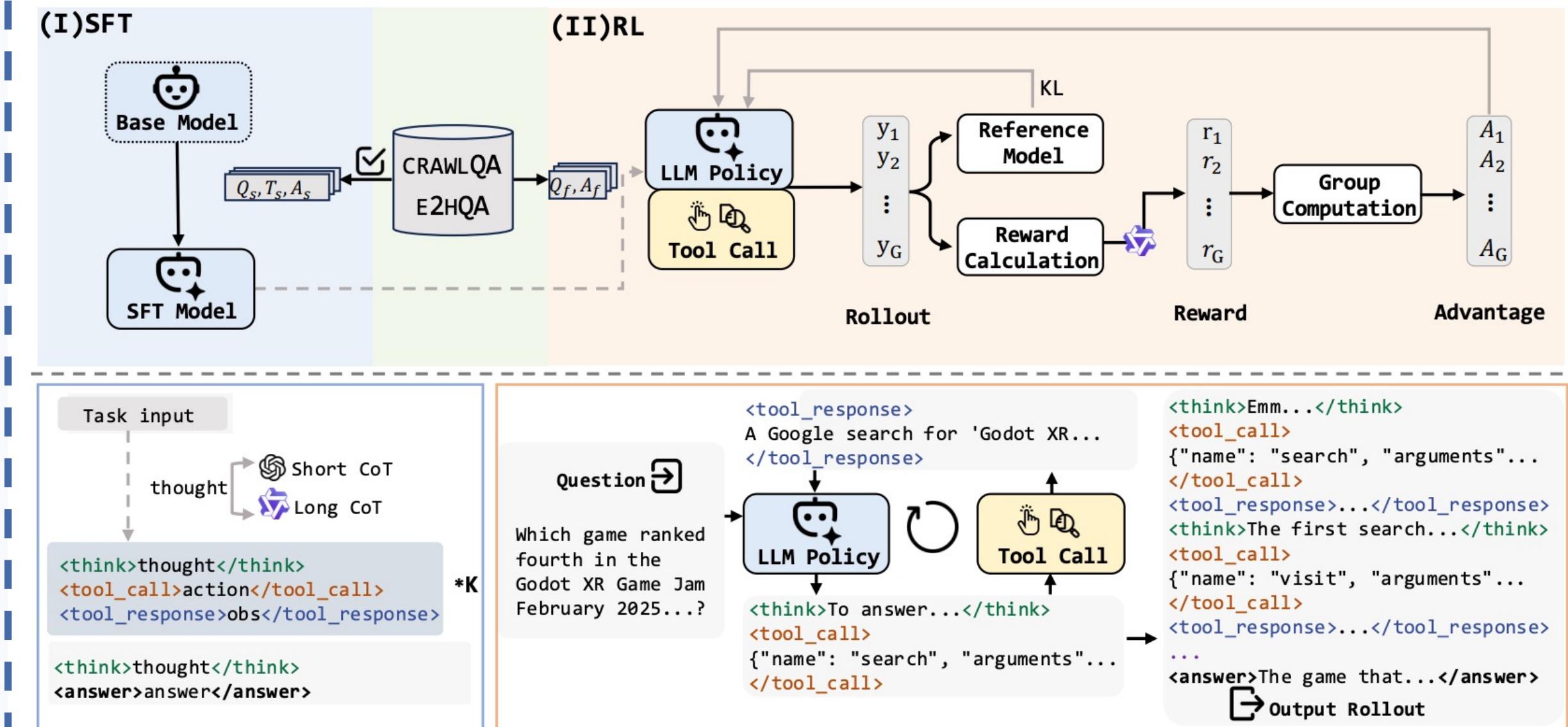


CRAWLQA We begin by collecting the root URLs of official and knowledgeable websites spanning *arxiv*, *github*, *wiki*, etc. To emulate human browsing behavior, we recursively navigate subpages by following accessible hyperlinks from each root site.

E2HQA We begin from large QA pairs in SimpleQA style where each answer is a concise, fact-seeking entity. By continuously searching, we can gradually rephrase an initially simple question into a more complex multi-step one.

High-quality synthetic data is the foundation of agent RL.

Training



(I) The SFT stage for cold start utilizes the reformatted **ReAct** datasets, where the thought includes both short and long CoT, respectively.

(II) The RL stage performs rollouts with the tool calls on the QA pairs that are not utilized during the SFT stage, and optimizes the policy using the DAPO algorithm..

Experiments

Our method significantly enhances agentic capabilities over the underlying base model, validating the strength and generality of our approach.

Backbone	Framework	GAIA			WebWalkerQA		
		Level 1	Level 2	Level 3	Avg.	Easy	Medium
Qwen-2.5-7B	Base	12.8	3.8	0.0	6.8	1.25	0.8
Qwen-2.5-32B	Base	20.5	9.6	8.3	13.6	3.8	2.5
Qwen-2.5-7B	RAG	12.8	11.8	8.3	11.8	23.1	14.3
Qwen-2.5-7B	Base	20.5	13.5	0.0	14.6	9.4	7.1
GPT-4o	Base	23.1	15.4	8.3	17.5	6.7	6.0
QwQ-32B	Base	30.8	15.4	25.0	22.3	7.5	2.1
DeepSeek-R1-671B	Base	43.6	26.9	8.3	31.1	5.0	11.8
OpenAI DR							
Qwen-2.5-7B	Search-01	23.1	17.3	0.0	17.5	-	-
Qwen-2.5-32B	Search-01	28.2	19.2	8.3	20.4	-	-
Qwen-2.5-32B	Search-01	33.3	25.0	0.0	28.2	-	-
QwQ-32B	Search-01	53.8	34.6	16.7	39.8	43.1	35.0
QwQ-32B	WebThinker-Base	53.8	44.2	16.7	44.7	47.2	41.1
QwQ-32B	WebThinker-RL	56.4	50.0	16.7	48.5	58.8	44.6
QwQ-32B	Simple DS	-	-	-	50.5	-	-
Open-sourced Agentic Frameworks							
Qwen-2.5-7B	R1-Searcher	23.1	17.3	0.0	17.5	-	-
Qwen-2.5-32B	Search-01	28.2	19.2	8.3	20.4	-	-
Qwen-2.5-32B	Search-01	33.3	25.0	0.0	28.2	-	-
QwQ-32B	Search-01	53.8	34.6	16.7	39.8	43.1	35.0
QwQ-32B	WebThinker-Base	53.8	44.2	16.7	44.7	47.2	41.1
QwQ-32B	WebThinker-RL	56.4	50.0	16.7	48.5	58.8	44.6
QwQ-32B	Simple DS	-	-	-	50.5	-	-
ReAct Agentic Frameworks							
Qwen-2.5-7B	Vanilla React	28.2	15.3	0.0	18.4	28.1	31.2
Qwen-2.5-7B	WebDancer	41.0	30.7	0.0	31.0	40.6	44.1
Qwen-2.5-32B	Vanilla React	46.1	26.9	0.0	31.0	35.6	38.7
Qwen-2.5-32B	WebDancer	46.1	44.2	8.3	40.7	44.3	46.7
QwQ-32B	Vanilla React	48.7	34.6	16.6	37.8	35.6	29.1
QwQ-32B	WebDancer	61.5	50.0	25.0	51.5	52.5	59.6
GPT-4o	Vanilla React	51.2	34.6	8.3	34.6	42.0	23.9

Results on BrowseComp (En.) and BrowseComp-zh (Zh.).

Performance across training steps using the DAPO algorithm.

End-to-end agentic RL is difficult, yet highly potential-rich.

Tongyi DeepResearch



Tongyi DeepResearch, developed by Tongyi Lab, is specifically designed for **long-horizon, deep information-seeking tasks**. Tongyi DeepResearch also has an extensive deep research agent family.

If you like our project, feel free to give us a on GitHub ☺