

More Efficient NLP and Agents

吴家隆(Jialong Wu)
wujialongml@gmail.com

Overall

- Training Stage
 - **Less Parameter:** LoRA
 - **Less Data:** Active Learning (SATR)
User Profiling (**PROPER**)
- Inference Stage
 - **Less Time:** Speculative Decoding (SEED)
 - **Less Memory:** KV-Cache Compression (**SCOPE**)
- Agent-Driven Autonomous Task Execution
 - **Less Human Intervention:**
 - **WebWalker**
 - **WebDancer**

Training Stage

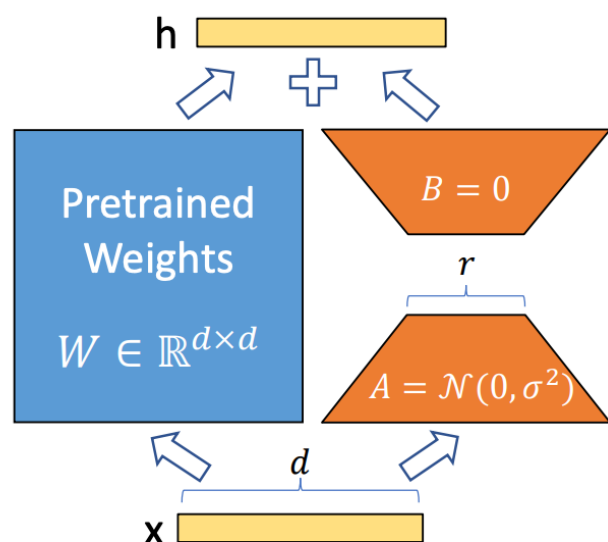


Figure 1: Low-rank adaptation.

- Finetuning Large Language Models (LLM) from scratch is quite resource-intensive, given the large number of parameters these models contain.
- Parameter efficient fine tuning is introduced.
- However, task-specific/user-specific fine-tuning also demands **high-quality data**.
- We propose **parameter and data efficient fine tuning**:
 - Active learning for task annotation
 - Progressive learning for personalized LLM

Training Stage

Active learning for task annotation

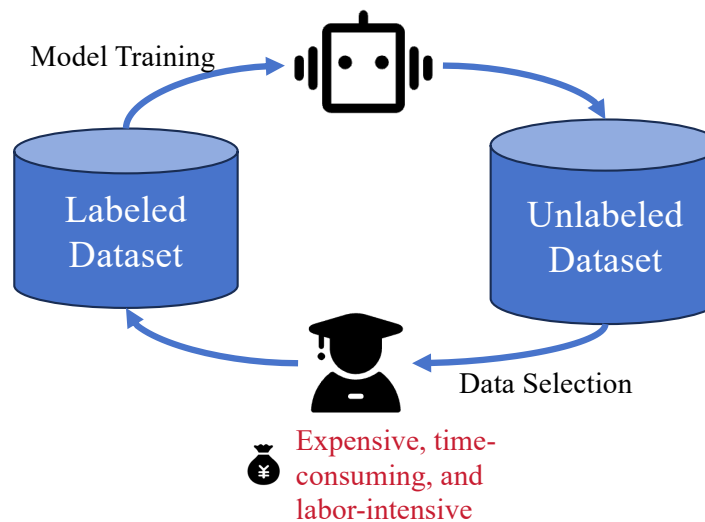


Figure 2: **Active Learning.**

Progressive learning for personalized LLM

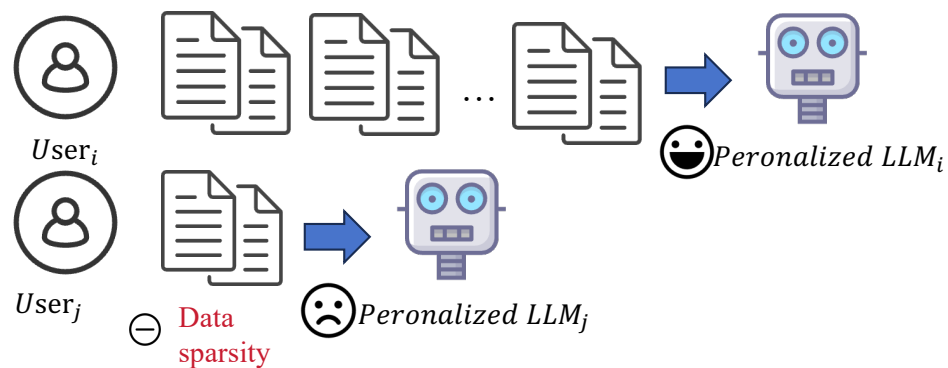


Figure 3: **Personalized LLM Training.**

Training Stage

STAR: Constraint LoRA with Dynamic Active Learning for Data-Efficient Fine-Tuning of Large Language Models

Probe experiments reveals:

- A clear **gap** between the base model and LoRA model
- Model **calibration** issue

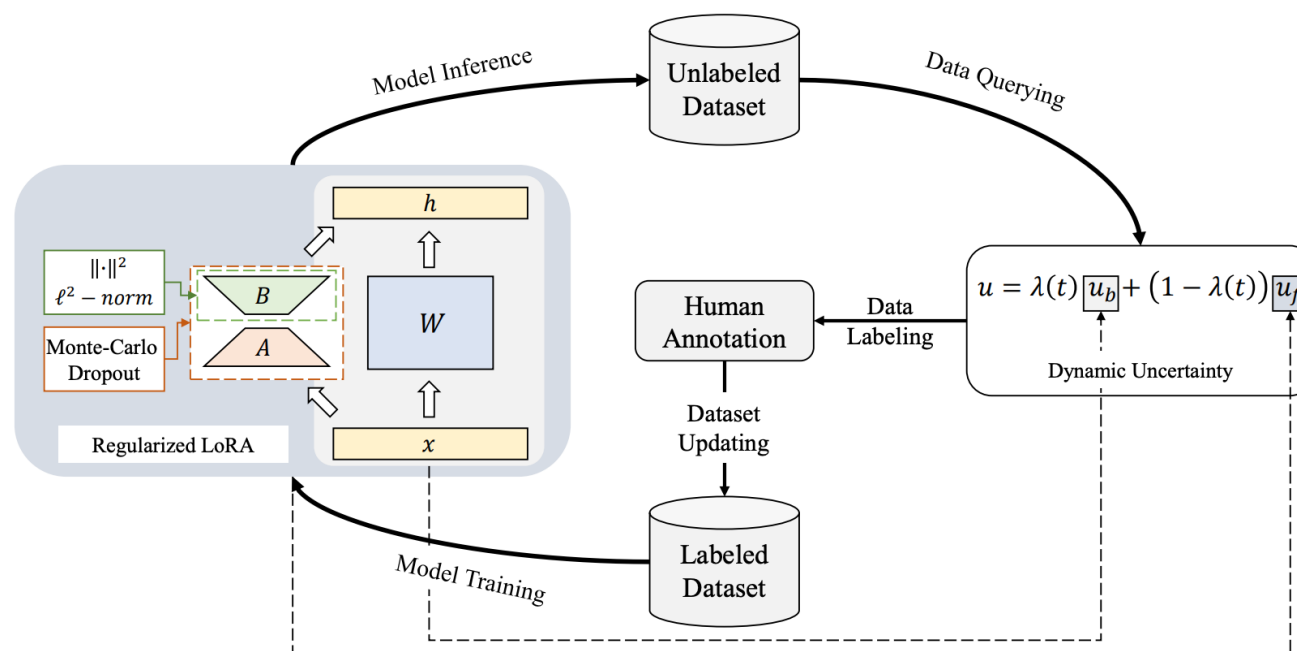


Figure 4: Overall of SATR.

Motivation: PROPER

- Current Large Language Models primarily offer a one-size-fits-all service.
- Personalized LLMs, tailoring the outputs to user-specific preferences, became a hot research topic.
- Two types of LLM personalization methods: (1) **prompt-based**, (2) **fine-tuning-based**.

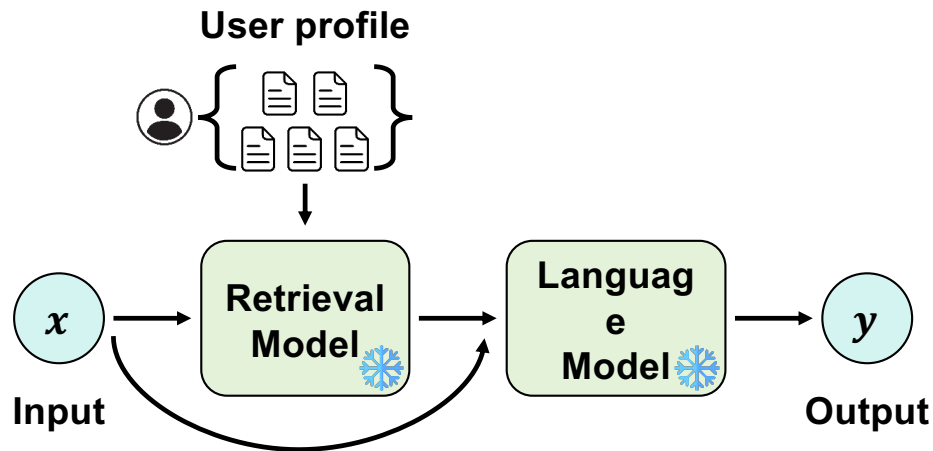


Fig. 1 (a) prompt-based personalized LLM,

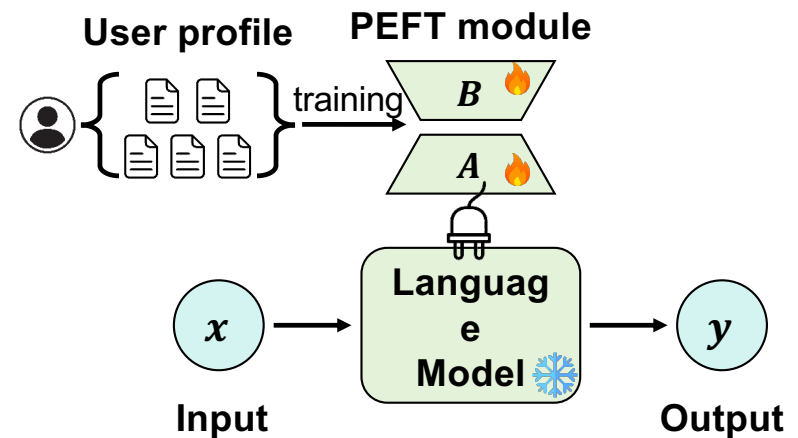


Fig. 1 (b) fine-tuning-based personalized LLM.

Fine-Tuning-based Personalized LLMs

- **Pros:** (1) data privacy, (2) better user behavior pattern generalization.
- **Cons:** (1) data scarcity for most of users (# avg tokens for task training is 20x of # avg tokens for individual users).
- **Solutions:** introduce meso-level LLMs (group-level LLMs) to bridge the macro-level LLMs (general LLMs) and micro-level LLMs (personalized LLMs).

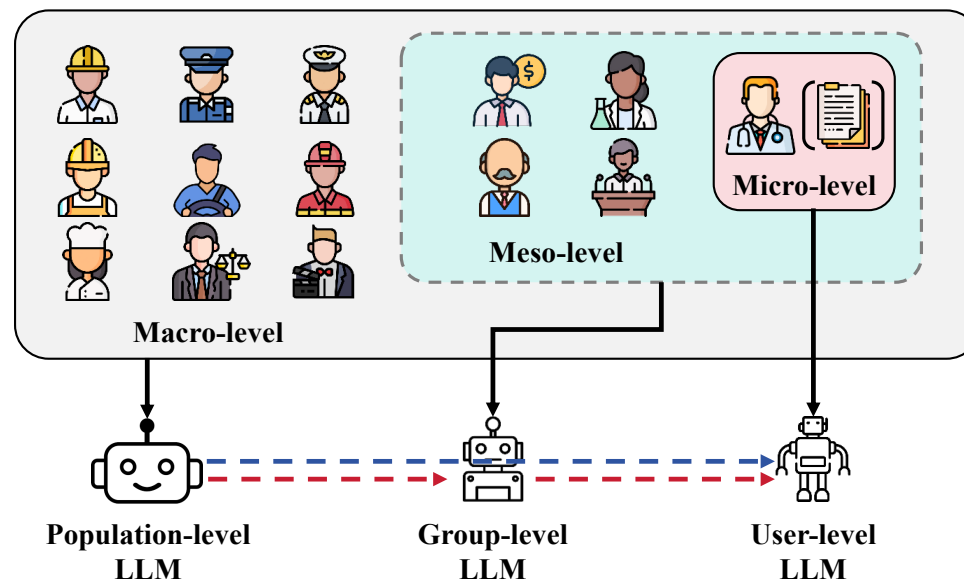


Fig 2. The comparison between different paradigms of LLM personalization.

Method

- A progressive learning framework, PROPER, consists of three stages: (1) population-level adaptation, (2) group-level adaptation, (3) user-level adaptation.
- Enable automatic user grouping via LoRAMoE and user-aware routers, while effectively integrating user and group-level knowledge through a LoRA-aware router.

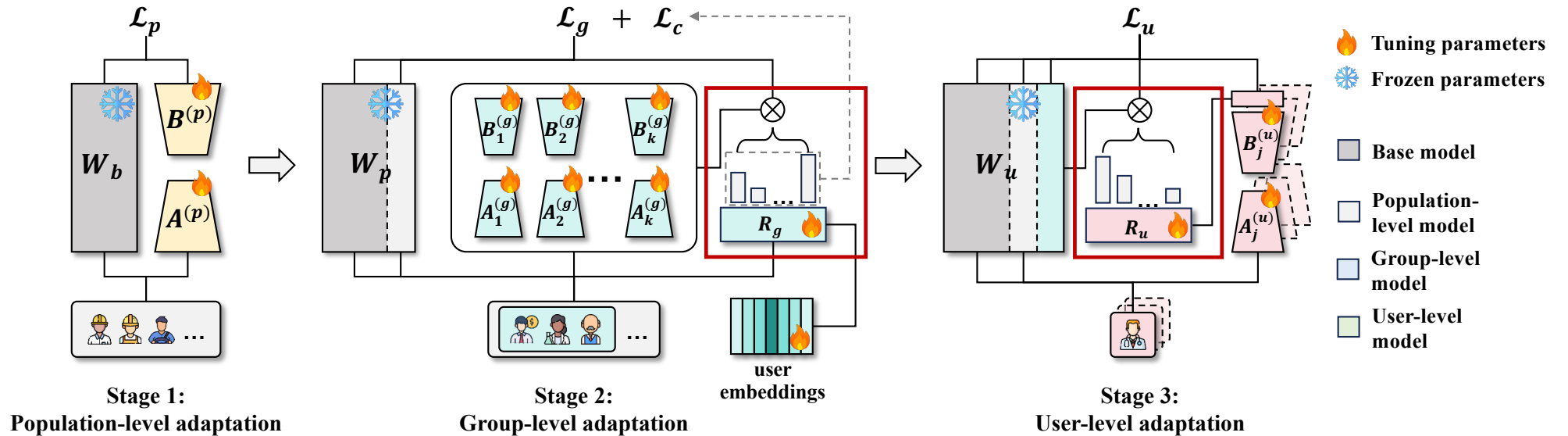


Fig 3. Overview of the training process of PROPER.

Method

Stage 1 (Population-Level Adaptation):

- The update process of the feed-forward network (FFN) block in a Transformer can be expressed as:

$$o = Wx = W_b + \Delta Wx$$

- In the population-level adaptation stage, parameter updates are formulated as:

$$o = W_b x + \frac{\alpha}{r} B^{(p)} A^{(p)} x$$

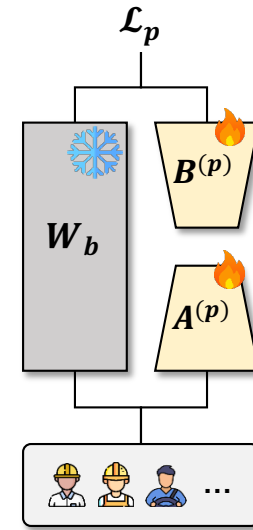
- The population-level LoRA is trained via fine-tuning using the cross-entropy loss:

$$\mathcal{L}_p = \sum_i \text{CE}\{\text{LLM}(q_i|\Omega_p), r_i\}$$

- With the loss, Stage 1 parameters are learned and merged into the backbone parameters for the next training stage:

$$\hat{B}^{(p)}, \hat{A}^{(p)} = \arg \min_{\Omega_p} \mathcal{L}_p$$

$$W_p = W_b + \hat{B}^{(p)} \hat{A}^{(p)}$$



Stage 1:
Population-level adaptation

Method

Stage 2 (Group-Level Adaptation):

- Employing LoRAMoE, represent each group with a LoRA experts:

$$o = W_p x + \sum_{i=1}^k \omega_i B_i^{(g)} A_i^{(g)} x$$

- Assign users to groups dynamically through a user-aware router:

$$\omega(x) = \text{softmax}(h),$$

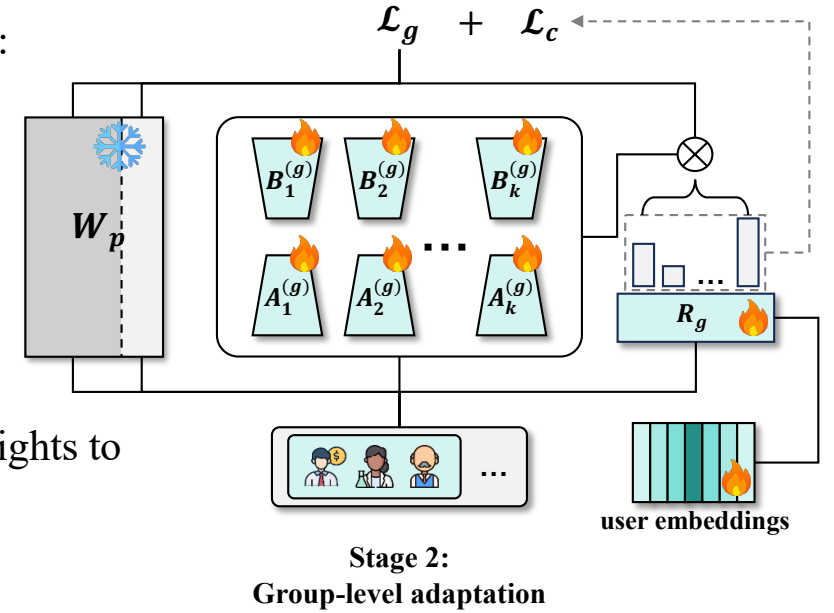
$$h = \text{softmax}(xM_g) + \text{softmax}(uM_u)$$

- Constraint loss to encourage the router to assign distinct expert weights to different users:

$$s_{(i,j)} = \omega_{u_i}^T \omega_{u_j} \quad \mathcal{L}_c = \sum_{i \neq j} |s_{(i,j)}|$$

- Learn Stage 2 parameters and merge into the backbone parameters with similar process in Stage 1:

$$\mathcal{L}_g = \sum_i \text{CE}\{\text{LLM}(q_i|\Omega_g), r_i\} \quad \hat{B}_j^{(g)}, \hat{A}_j^{(g)} = \arg \min_{\Omega_g} \mathcal{L}_g \quad W_g = W_p + \frac{a}{r} \sum_{j=1}^k \omega_j B_j^{(g)} A_j^{(g)}$$



Method

Stage 3 (User-Level Adaptation):

- Assign a unique LoRA to each user:

$$o = W_g x + B_j^{(u)} A_j^{(u)} x$$

- A new LoRA-aware router that dynamically integrates group-level LoRAs and user-level LoRAs:

$$\beta_u(x) = \text{softmax}(W_l h_u)$$

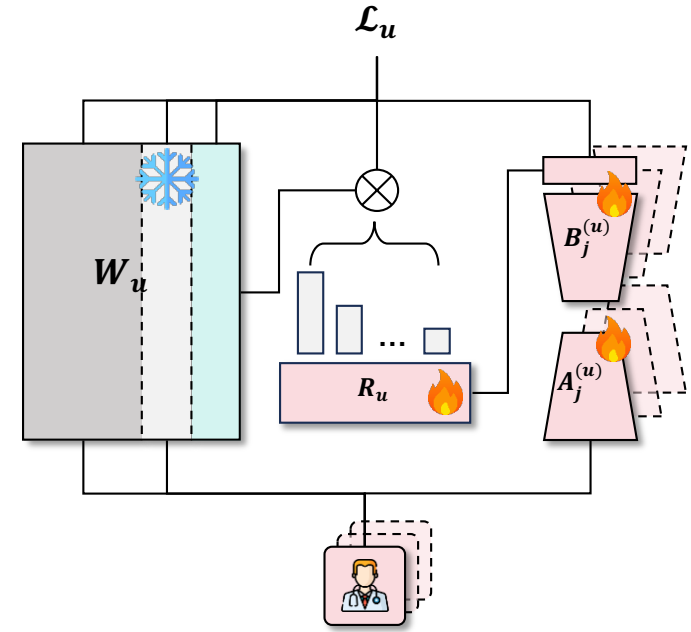
$$h_u = \text{LoRA}_u(x),$$

- Learn Stage 3 parameters and merge into the backbone parameters with similar process in Stage 1&2:

$$\mathcal{L}_p^{(j)} = \sum_i \text{CE}\{\text{LLM}(q_i^{(j)} | \Omega_p^{(j)}), r_i^{(j)}\},$$

$$\hat{B}_j^{(u)}, \hat{A}_j^{(u)} = \arg \min_{\Omega_p^{(j)}} \mathcal{L}_p^{(j)},$$

$$W_u^{(j)} = W_g + B_j^{(u)} A_j^{(u)} + \sum_{m=1}^k \beta_m B_m^{(g)} A_m^{(g)},$$



Stage 3:
User-level adaptation

Experiments

Evaluation Dataset:

- LaMP benchmark

Baselines:

- Prompt-based
 - In-Context-Learning (ICL)
 - Retrieval-Augmented Generation (RAG)
 - Profile-Augmented Generation (PAG)
- Fine-tuning-based
 - OPPU (kv, mlp)
 - PROPER (Stage 1, 2, 3)

Input	Output	Profile
LaMP-1: Personalized Citation Identification For an author who has written the paper with the title "[TITLE]", which reference is related? Just answer with [1] or [2] without explanation. [1]: "[REF1]" [2]: "[REF2]"	[1]	title: [TITLE] abstract: [ABSTRACT]
LaMP-2: Personalized Movie Tagging Which tag does this movie relate to among the following tags? Just answer with the tag name without further explanation. tags: [sci-fi, based on a book, comedy, action, twist ending, dystopia, dark comedy, classic, ...] description: [MOVIE]	comedy	description: [MOVIE] tag: [TAG]
LaMP-3: Personalized Product Rating What is the score of the following review on a scale of 1 to 5? Just answer with 1, 2, 3, 4, or 5 without further explanation. review: [REVIEW]	4	text: [REVIEW] score: [SCORE]
LaMP-4: Personalized News Headline Generation Generate a headline for the following article: [ARTICLE]	The Best Cheap Wine: Two Buck Chuck vs Three Wishes	title: [TITLE] text: [ARTICLE]
LaMP-5: Personalized Scholarly Title Generation Generate a title for the following abstract of a paper: [ABSTRACT]	Attention is All You Need	title: [TITLE] abstract: [ABSTRACT]
LaMP-7: Personalized Tweet Paraphrasing Paraphrase the following tweet without any explanation before or after it: [TWEET]	I hope so! what time do you get out? I get out at 335	text: [TWEET]

Fig 5. LaMP input & output examples.

Inference Stage

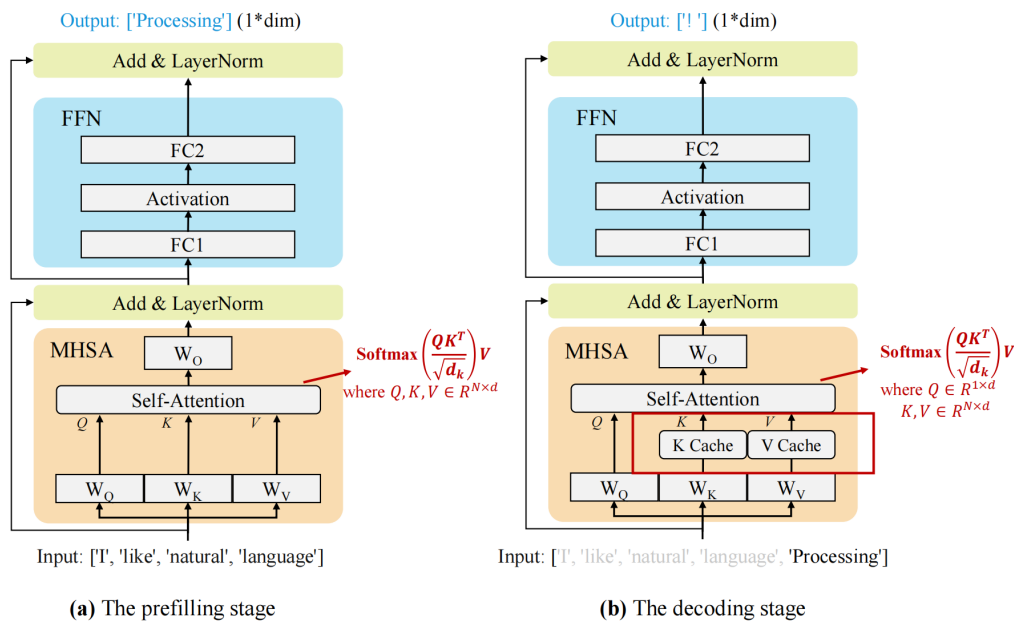


Figure 6: **Autoregressive** token-by-token generation leveraging the **key-value (KV) cache**.

The token-by-token inference of large models results in slower processing speeds, as each step of generation requires the use of the KV-cache.

Multi-output: ToT require independently traversing each branch. This leads to an exponential increase in inference time due to the autoregressive nature of the model.

Long-output: Long-form text leverage the KV-cache to accelerate attention computations, results in a sharp increase in memory usage.

Inference Stage

Multi output

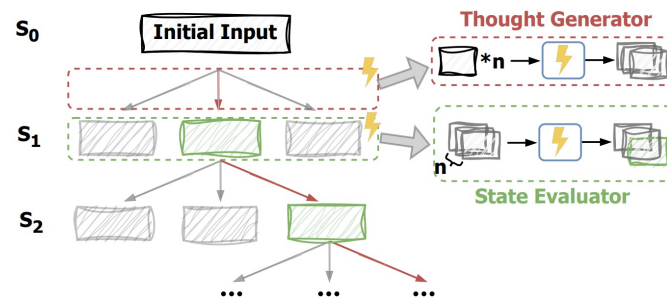
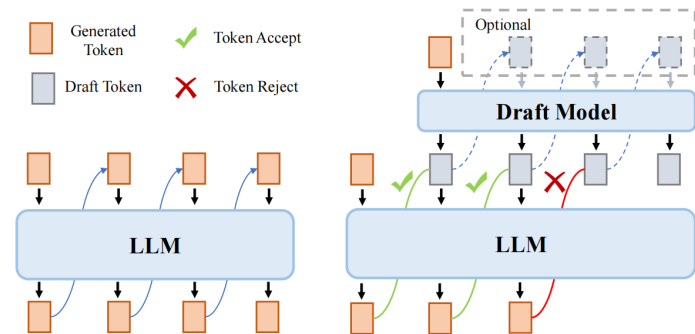


Figure 7: Tree of thoughts.

Figure 8: Paradigm of **speculative decoding**.

Long output

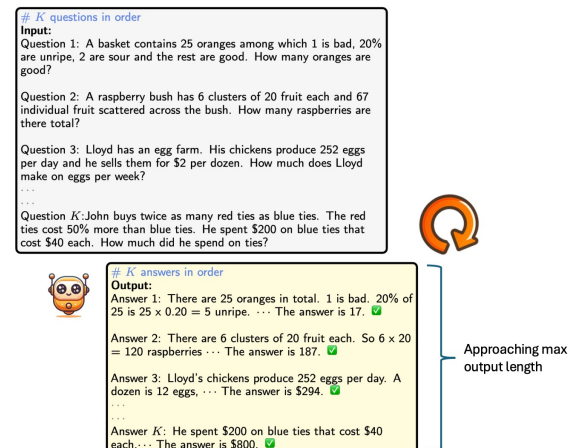


Figure 9: Long-context Generation.

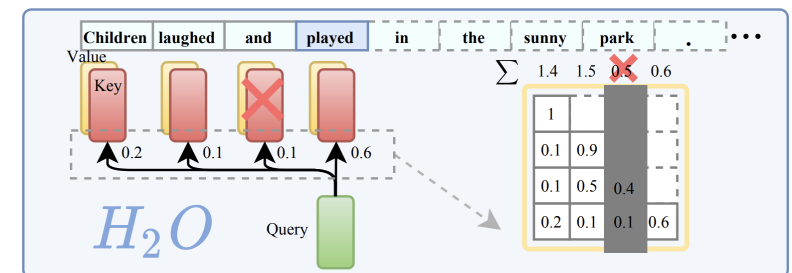


Figure 10: KV cache compression (eviction) guided by attention scores.

Inference Stage

SEED: Accelerating Reasoning Tree Construction via Scheduled

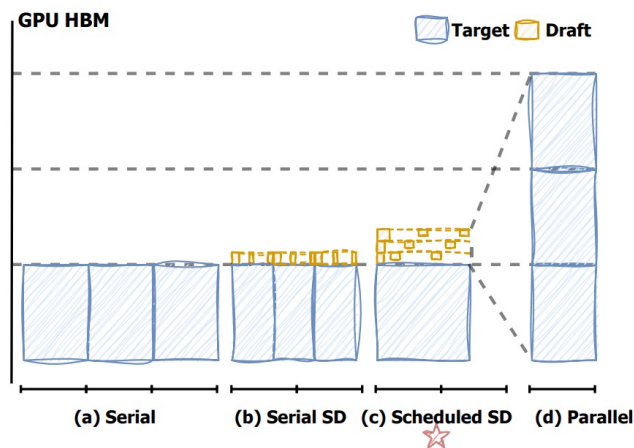


Figure 11: Illustration of four LLM execution strategies for generating 3 sequences in Reasoning Tree construction.

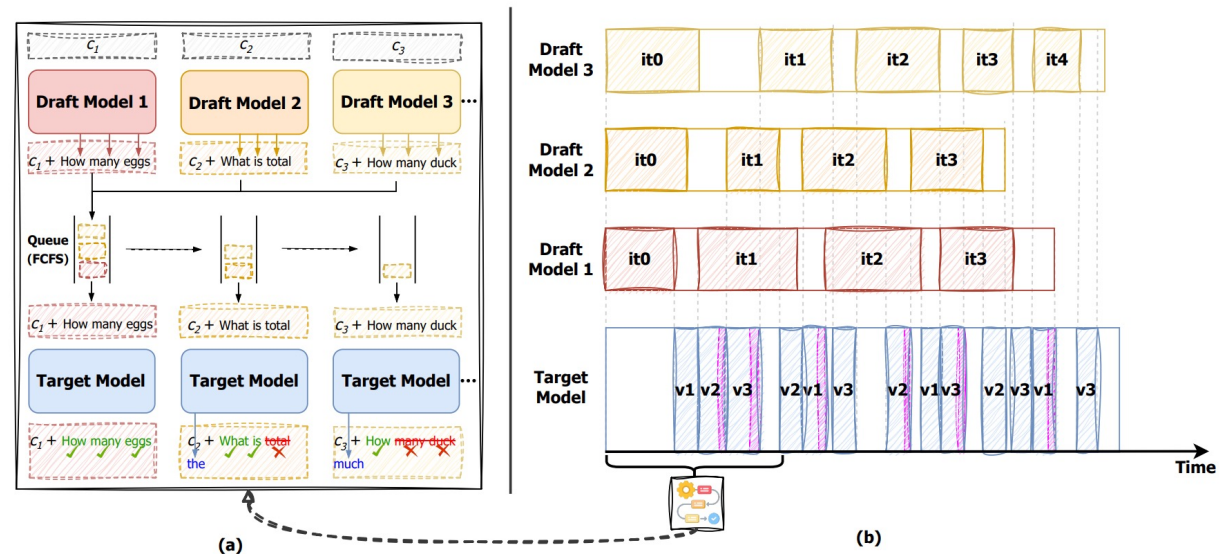


Figure 12: Overall scheduling diagram of **SEED**.

Motivation

- When Large Language Models infer on long-context tasks, the Key-Value (KV) cache occupies a larger amount of GPU memory and becomes a substantial bottleneck.
- Previous methods fall into two categories:
 - (1) The **Prefill-Only** compression method ➡ memory pressure for long outputs
 - (2) The **Unified** compression method ➡ fine-grained content eviction

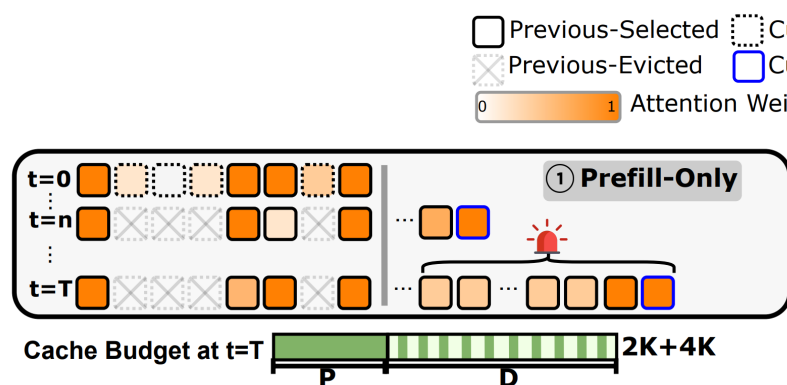


Fig. 1 (a) The Prefill-Only Compression method,

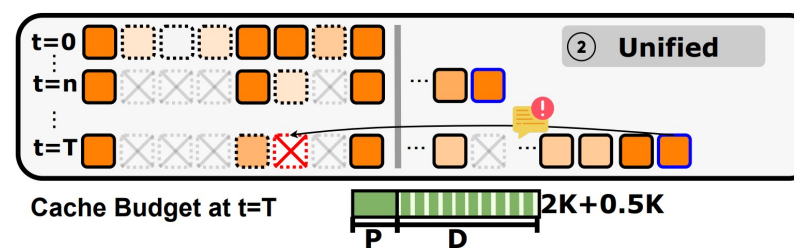


Fig. 1 (b) The Unified Compression method

Separating the Prefill and Decoding Phases

- **Prefill Phase:** Efficiently preserves essential information in the KV cache during the prefill.
- **Decoding Phase:** Enables optimized allocation of KV cache generated during decoding.
- **Solutions:** SCOPE, a simple yet efficient framework that Separately performs KV Cache Optimization during the Prefill and dEcoding phases.

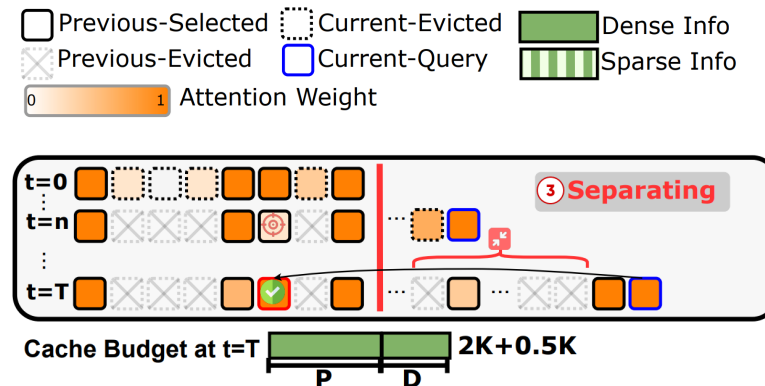


Fig. 1 (c) Separating the prefill and decoding phases

Pilot Observation: KV Cache in Inference Perspective

Prefill Stage

The **20%** compression rate during the prefill phase resulted in nearly **95% degradation** in accuracy on the GSM8k+ task within LONGGENBENCH.



Observations (i): For tasks that require specific fine-grained context, such as reasoning tasks, excessive compression during the prefill phase significantly compromises performance.



Fig. 2 (a) Performances across various compression ratios during the prefill phase on three tasks under the full decoding cache condition.

Pilot Observation: KV Cache in Inference Perspective

Decoding Stage

Across all three layers, the retained heavy hitters predominantly originate from the KV cache generated **during the decoding phase**.



Observations (ii): During the decoding phase of long text generation, the use of the greedy algorithm may lead to a deviation in heavy hitters.

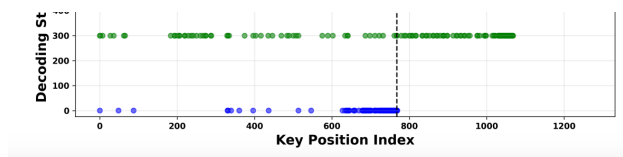


Fig. 2 (a) Position distribution of the heavy hitters, selected by top 15% attention scores, at decoding steps 1, 300, and 500 across layers 0, 13, and 31.

Pilot Observation: KV Cache Budget Reallocation

Observations (i): For tasks that require specific fine-grained context, such as reasoning tasks, excessive compression during the prefill phase significantly compromises performance.

Observations (ii): During the decoding phase of long text generation, the use of the greedy algorithm may lead to a deviation in heavy hitters.

Pilot Observation: KV Cache Budget Reallocation

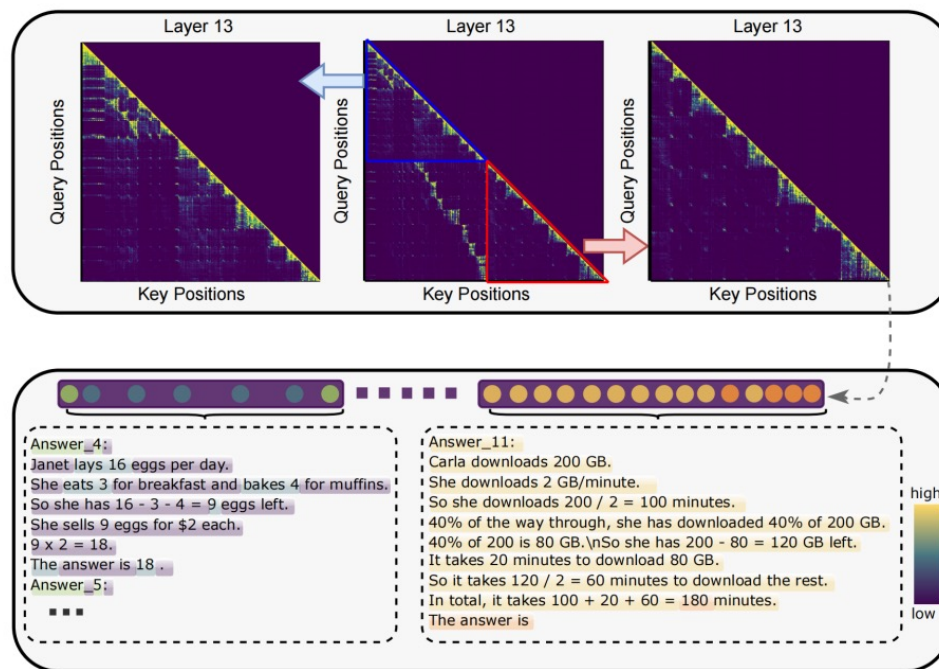
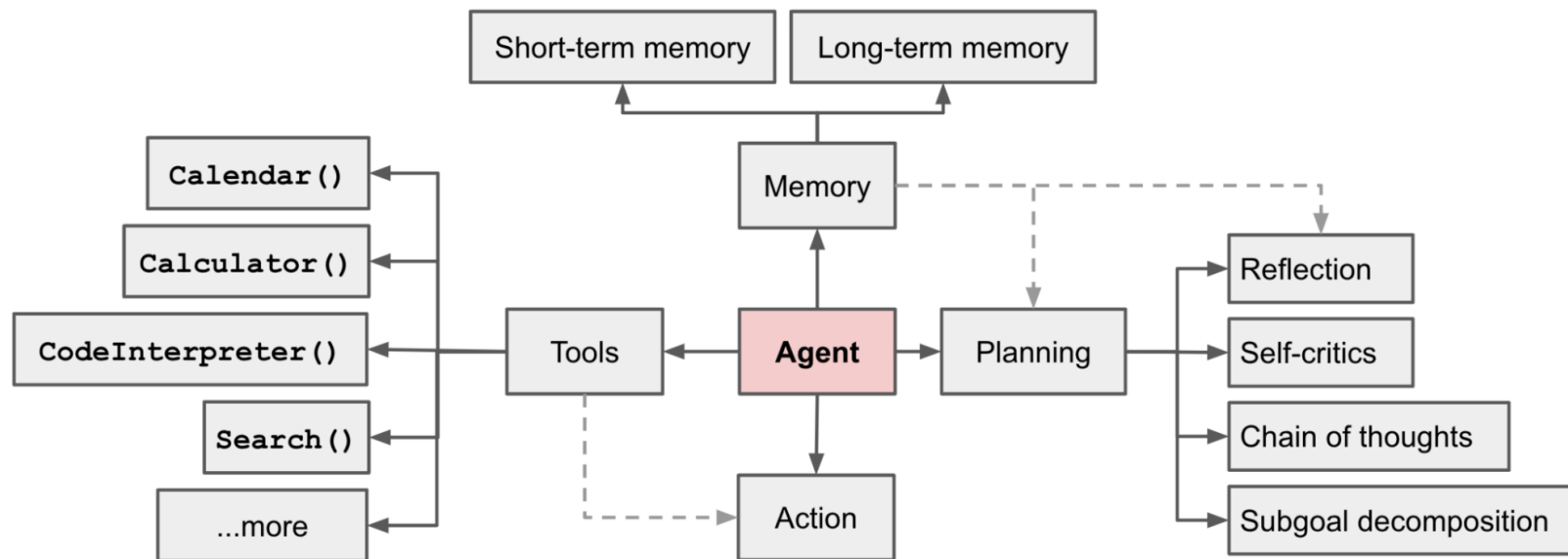


Fig. 2 (c) Attention heatmaps for layer 13 of a GSM8k+ sample in LONGGENBENCH and details of the correspondence between attention scores and generated token positions.

Agent



Agent

Agents: An Open-source Framework for Autonomous Language Agents

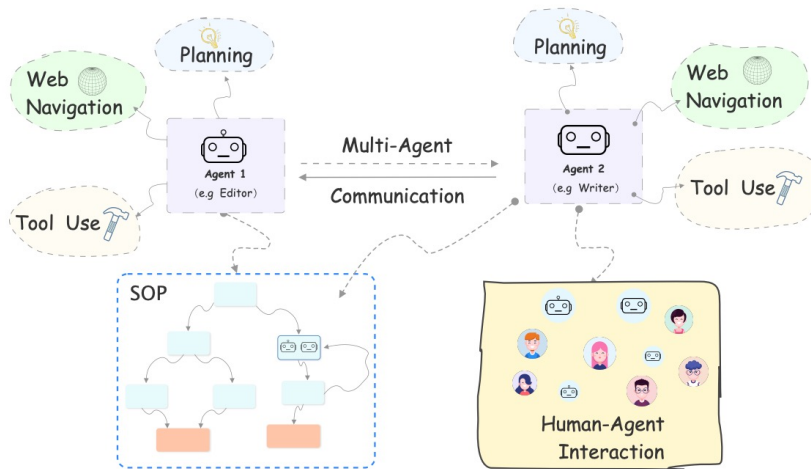


Figure 16: Components in agents framework.

- **Long-short Term Memory:** Long-term memory implemented via VectorDB + Semantic Search and short-term memory (working memory) maintained and updated by an LLM.
- **Tool Usage:** Use any external tools via function-calling.
- **Web Navigation:** Use search engines to navigate the web and get useful information.
- **Multi-agent Communication**
- **Human-Agent interaction**
- **Symbolic Control:** SOP (Standard Operation Process) that defines subgoals/subtasks for the overall task to customize fine-grained workflows for the language agents.

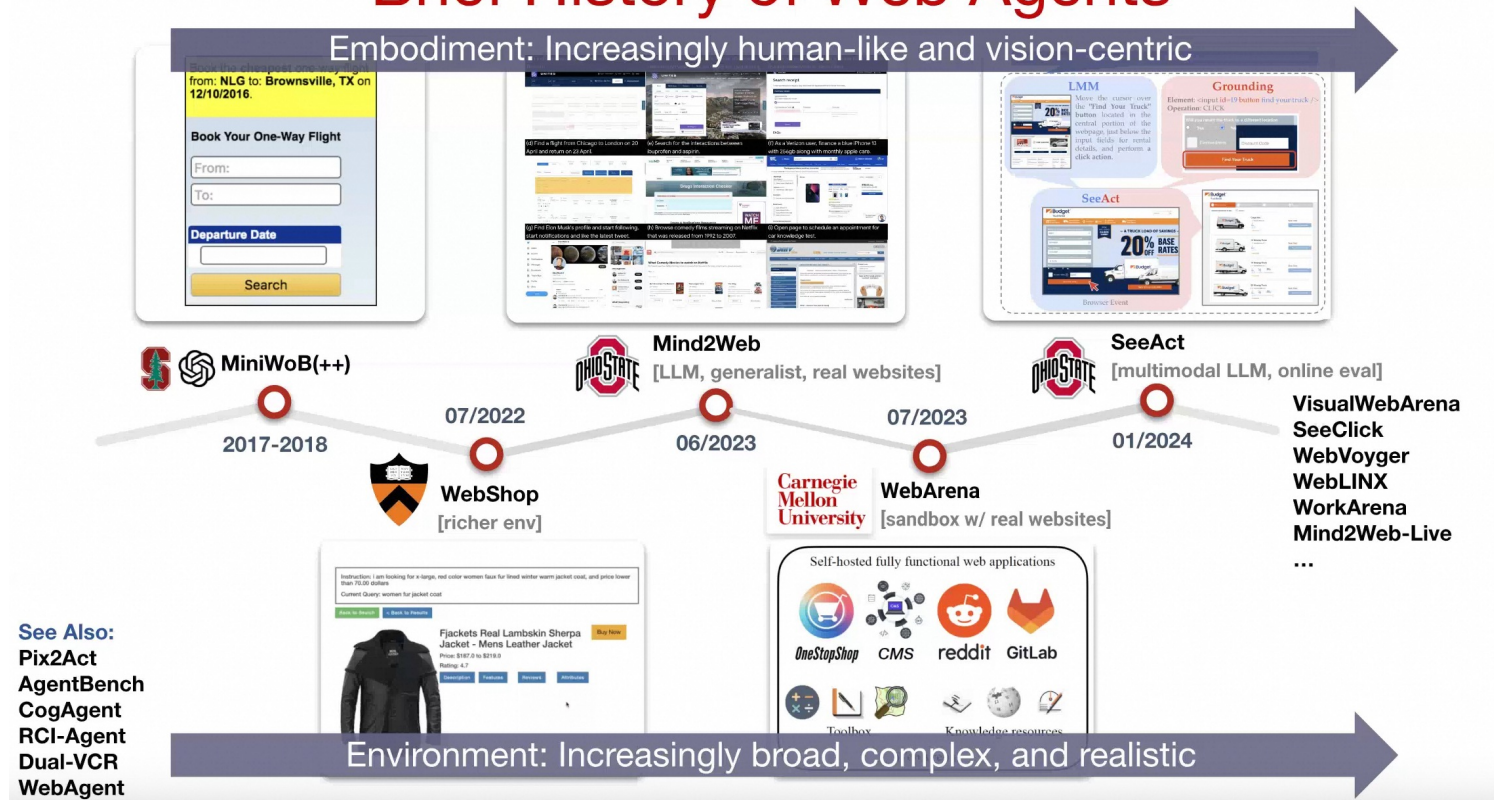
Agent

Brief History of Web Agents



Agent

Brief History of Web Agents

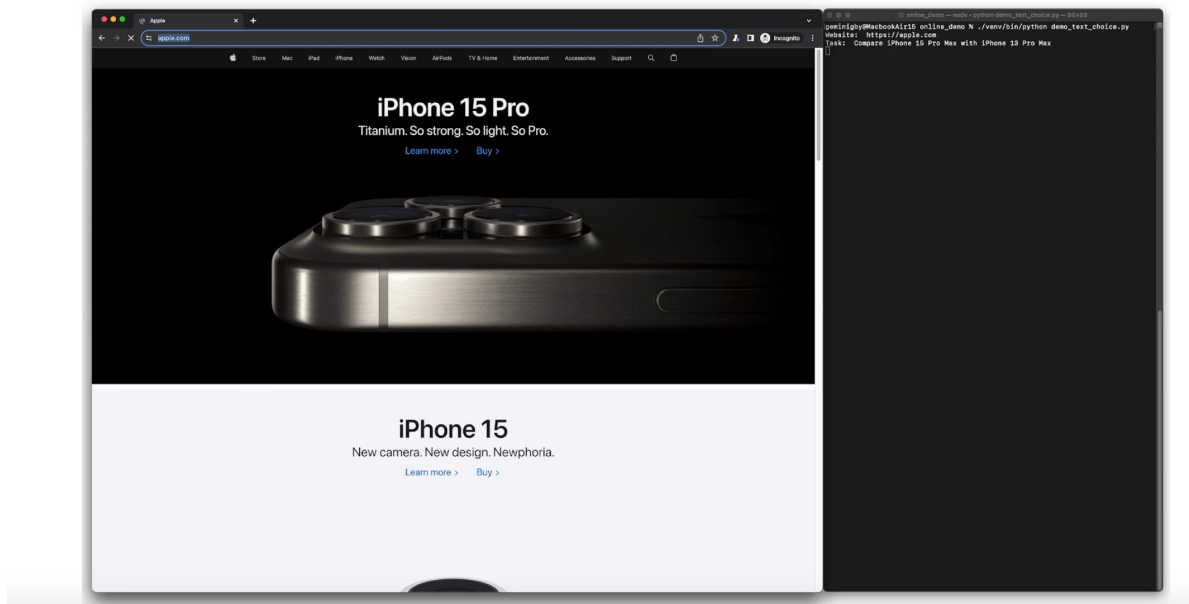


Agent

Generalist Web Agents: Mind2Web & SeeAct

(NeurIPS'23)

(ICML'24)



Website: <https://apple.com>

Task: Compare iPhone 15 Pro Max with iPhone 13 Pro Max

Does gpt-3.5-turbo support structured outputs, like response_format: {type: "json_schema", ...}?



...Yes, GPT-3.5-turbo supports structured outputs.



We recommend always using Structured Outputs instead of JSON mode when possible.

However, Structured Outputs with `response_format: {type: "json_schema", ...}` is only supported with the `gpt-4o-mini`, `gpt-4o-mini-2024-07-18`, and `gpt-4o-2024-08-06` model snapshots and later.

STRUCTURED OUTPUTS		JSON MODE
Outputs valid JSON	Yes	Yes
Adheres to schema	Yes (see supported schemas)	No
Compatible models	<code>gpt-4o-mini</code> , <code>gpt-4o-2024-08-06</code> , and later	<code>gpt-3.5-turbo</code> , <code>gpt-4o</code> and <code>gpt-4o-mini</code>
Enabling	<code>response_format: { type: "json_schema", json_schema: {"strict": true, "schema": ...} }</code>	<code>response_format: { type: "json_object" }</code>

What is the latest publication written by openai?



...OpenAI's latest research paper is titled "PaperBench: Evaluating AI's Ability to Replicate AI Research," published on April 2, 2025.



Research

All Publication Conclusion Milestone Release

Filter Sort

Publication

Apr 16, 2025

OpenAI o3 and o4-mini System Card

OpenAI o3 and OpenAI o4-mini combine state-of-the-art reasoning with full tool capabilities—web browsing, Python, image and file analysis, image generation, canvas, automations, file search, an...

Publication

Apr 15, 2025

Our updated Preparedness Framework

Sharing our updated framework for measuring and protecting against severe harm from frontier AI capabilities.

- Can not find information hidden within deep webpages or the latest updates.

Agent

How can an agent navigate across web pages to seek deep, non-obvious information?

- Unlock the capability of persistent, multi-hop web exploration




Motivation

Key challenge in RAG:

Traditional online search may not trace the
Deeper content embedded within website.

When is the paper submission deadline for the ACL 2025 Industry Track, and what is the venue address for the conference?


 <https://2025.aclweb.org/>



Motivation

How to solve it:

Interacting with the web pages and **digging through** them can effectively address **deep information seeking**.

We constrain actions to click  to evaluate the agent's navigation and information-seeking capabilities.

- We propose **Web Traversal task**.
- We construct a challenging benchmark, **WebWalkerQA**.
- To tackle the challenge of web-navigation tasks requiring long context, we propose **WebWalker**.



Datasets

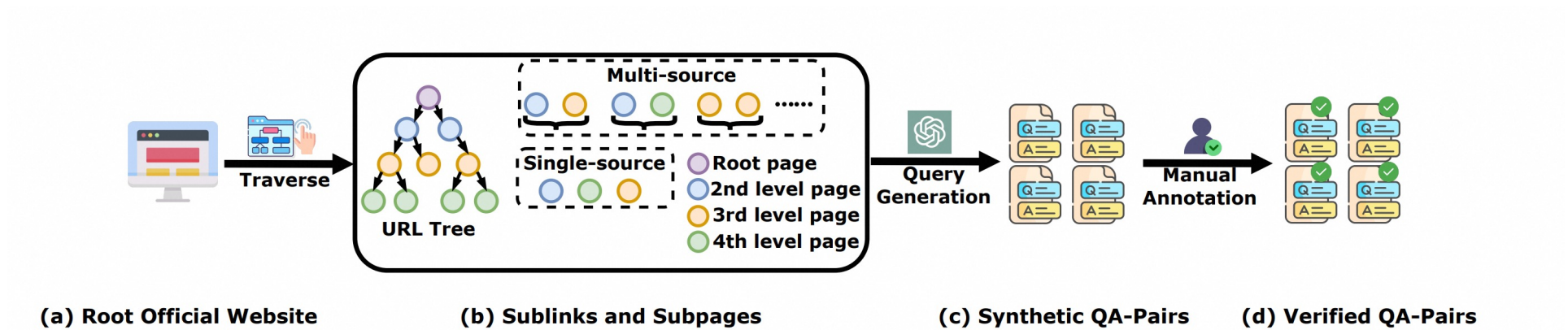
WebWalkerQA



	Language	Format	Depth	Width	Hop	# Pages
Mind2Web (Deng et al., 2023)	En	Multi-choice	✗	✗	✗	100
WebArena (Zhou et al., 2024a)	En	Action	✗	✗	✗	6
AssistantBench (Yoran et al., 2024)	En	QA	✗	✓	✓	525
MMinA (Zhang et al., 2024c)	En	Action	✗	✓	✓	100
GAIA (Mialon et al., 2024)	En	QA	✗	✓	✓	-
WebWalkerQA	En&Zh	QA	✓	✓	✓	1373

Comparison between WebWalkerQA and other benchmarks.

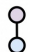

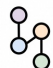



Datasets



Data Generation Pipeline for WebWalkerQA.

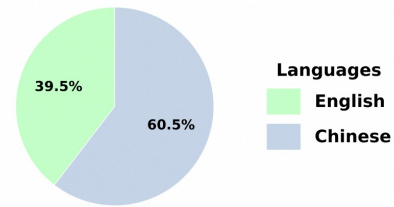
Datasets



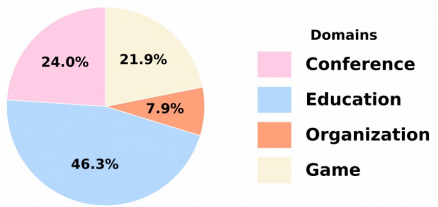
Single-source QAs			Multi-source QAs		
					
Easy	Medium	Hard	Easy	Medium	Hard
80	140	120	80	140	120

Dataset statistics on difficulty level.

Language Distribution



Domain Distribution



Language and domain **distribution**.

Datasets


Web Traversal Task:

Given an initial website URL and a query Q , which needs to be answered by exploring the website. The goal of this task is to gather enough information through page traversal to ultimately answer the query Q .


Evaluation:

Correctness -> acc. Evaluated by GPT-4o

Efficiency -> Action count of successful agentic executions



When is the paper [submission deadline for the ACL 2025 Industry Track](#), and what is the [venue address for the conference](#)?

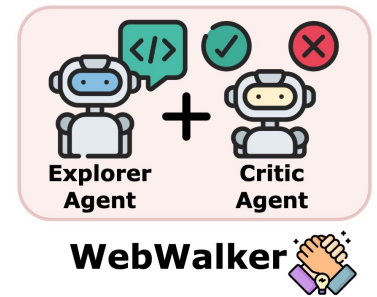


<https://2025.aclweb.org/>

The paper submission deadline for the ACL 2025 Industry Track is [March 21, 2025](#) and the conference will be held in [Bruno-Kreisky-Platz 1](#).

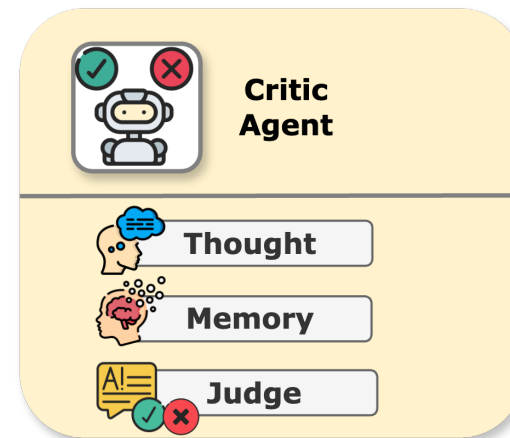
Methods

WebWalker: a multi-agent framework



Think then Explore

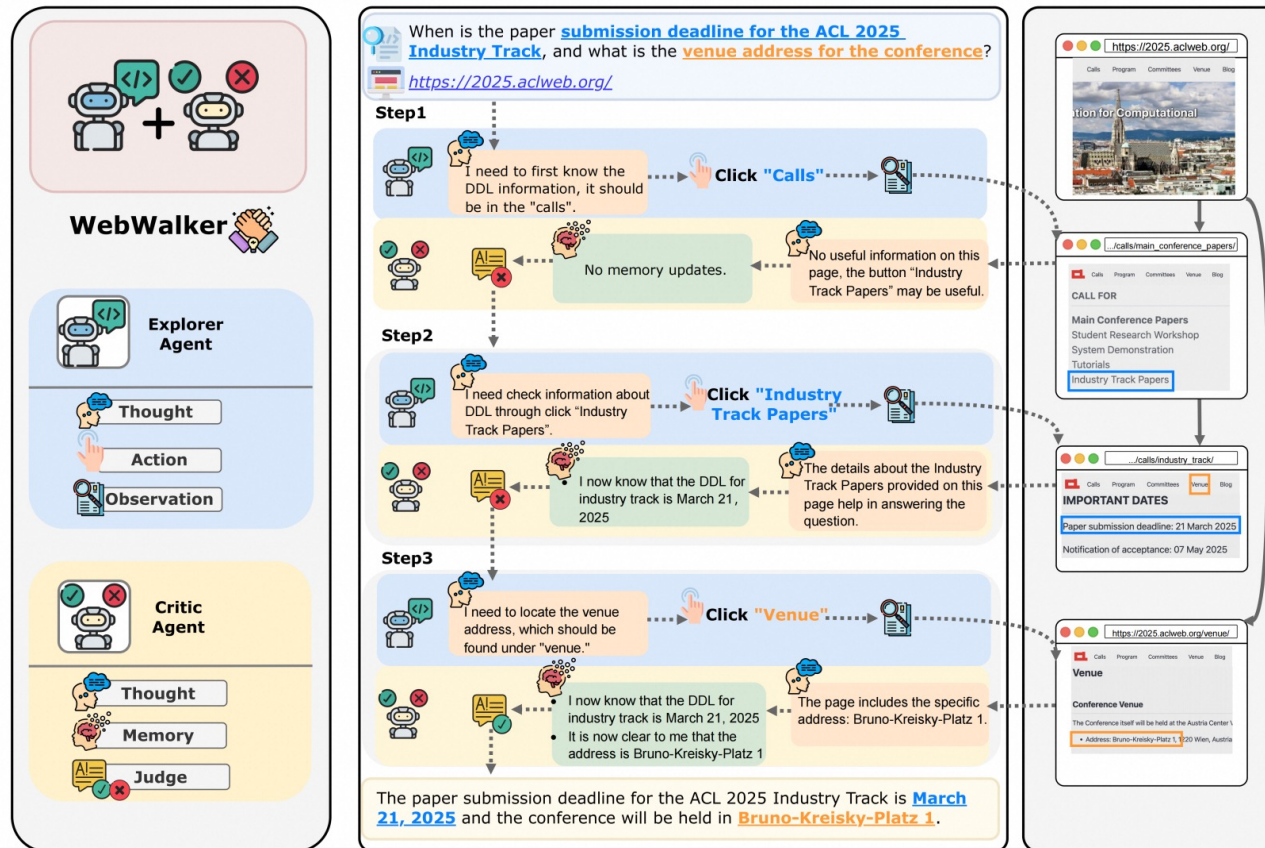
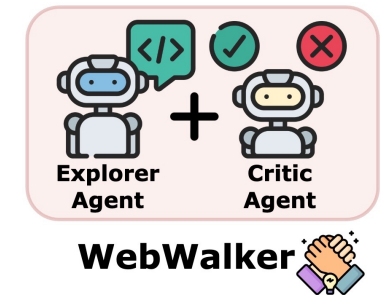
ReAct format



Think then Critique

Motivated by pair programming

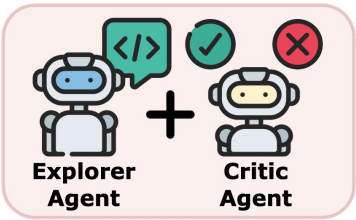
Methods




The explorer agent traverses the web pages in **Thought-Action-Observation** (T, A, O) paradigms.


The critic agent **updates the memory** until sufficient information is accumulated to effectively **address the query**.

Methods



WebWalker

 **WebWalker**

 Memory

No Memory

 Website

<https://2025.aclweb.org/>

 Query

When is the Industry Track paper submission deadline for ACL 2025, and what is the venue address?

Start!!!!

Agent

February 2, 2025 Release

Introducing deep research

Try on ChatGPT ↗

Agent

ChatGPT 4o

Share

MI

Activity 3 Sources

- Supply chain challenges
- Emerging technologies (AI, automation, etc.)
- Impact of inflation and economic changes
- Sustainability and ethical retailing Also, do you need a global perspective, or should the focus be on a specific region?

All, global perspective

I'll compile a detailed research report on how the retail industry has transformed over the past three years on a global scale. This will cover:

- The rise of e-commerce and shifts in brick-and-mortar retail
- Changes in consumer behavior and expectations
- Supply chain disruptions and adaptations
- The role of AI, automation, and emerging technologies
- The impact of inflation and economic fluctuations
- Sustainability and ethical retailing trends

I'll structure the report with bullet points and tables where necessary to enhance clarity. I'll let you know once it's ready for review.

Searched for retail foot traffic 2022 ... 3 sources

Message ChatGPT

Attach Search Deep research

ChatGPT can make mistakes. Check important info.

In 2023, e-commerce reached a 10.7% share of global retail sales, which is projected to increase to 22.6% by 2027.

Considering potential access issues.

Q Searched for Statista e-commerce share data

I reviewed e-commerce share data for the US from 2019 to 2023, noting a peak around 2020 and stabilization in subsequent years. Also, I looked into global e-commerce trends, particularly China's significant influence.

Read oberlo.com

OpenAI consistently upholds a dedication to ethical research and the well-being of all individuals.

Q Searched for global e-commerce share 2019 vs 2020

Checking for trends in e-commerce share growth from 2019 to 2020.

Read trade.gov

I'm thinking through the impact of global e-commerce growth on brick-and-mortar stores, noting trends like store closures and shifts towards omnichannel strategies, experiential retail, and smaller stores.

Q Searched for retail foot traffic 2022 rebound

Reading...

Deep Research is currently in Alpha.

Agent

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

- unlock the autonomous multi-turn information seeking agency



WebDancer

Pipelines

Step I: Construct diverse and challenging deep information seeking QA pairs based on the real-world web environment;

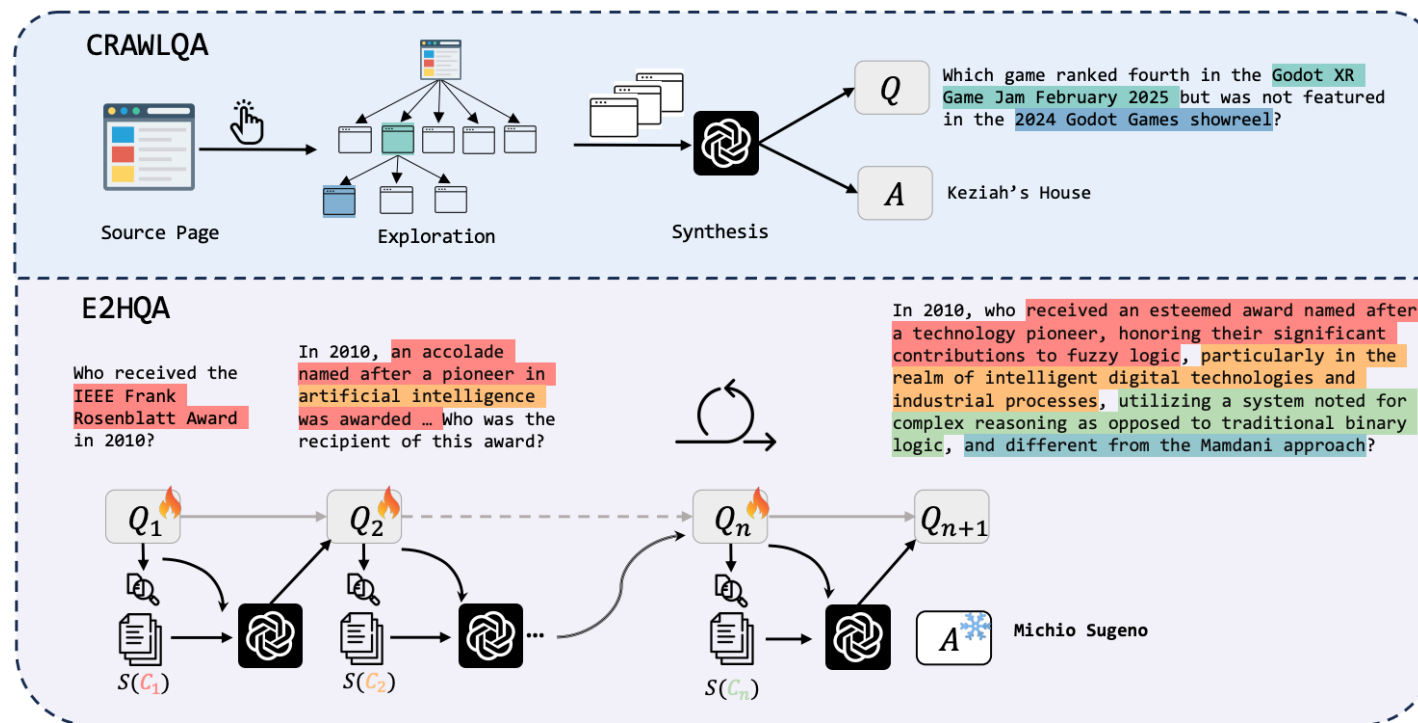
Step II: Sample high-quality trajectories from QA pairs using both LLMs and LRMs to guide the agency learning process;

Step III: Perform fine-tuning to adapt the format instruction following to agentic tasks and environments;

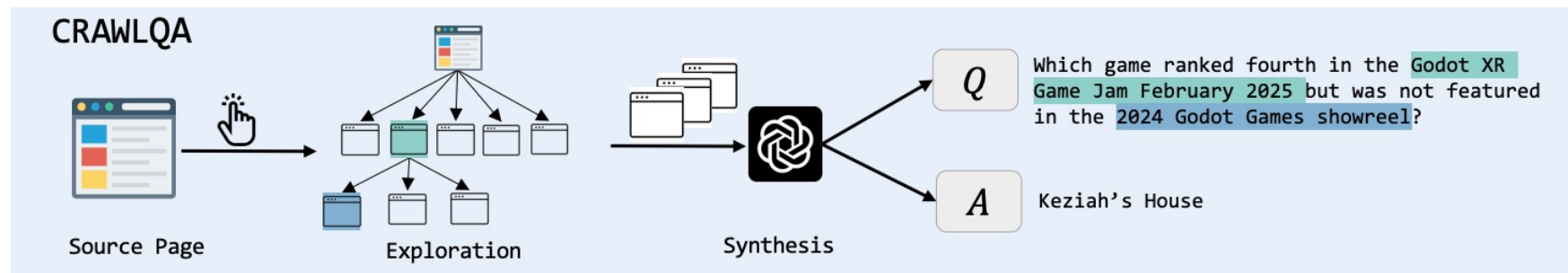
Step IV: Apply RL to optimize the agent's decision-making and generalization capabilities in real-world web environments

Datasets

Previous training datasets are relatively simple and do not capture the real-world challenges.

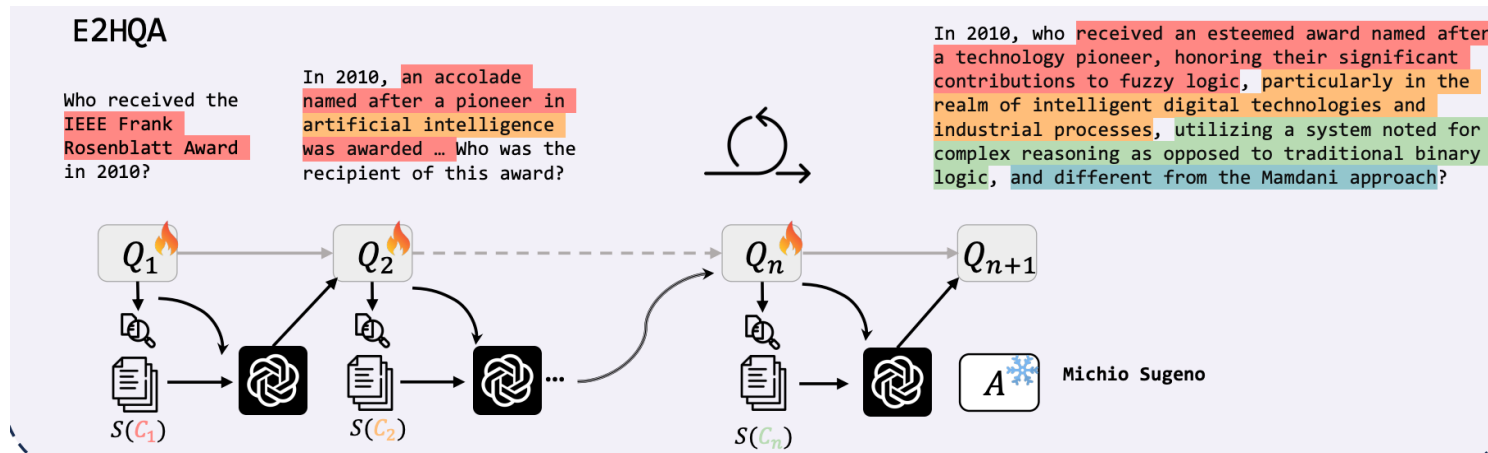


Datasets



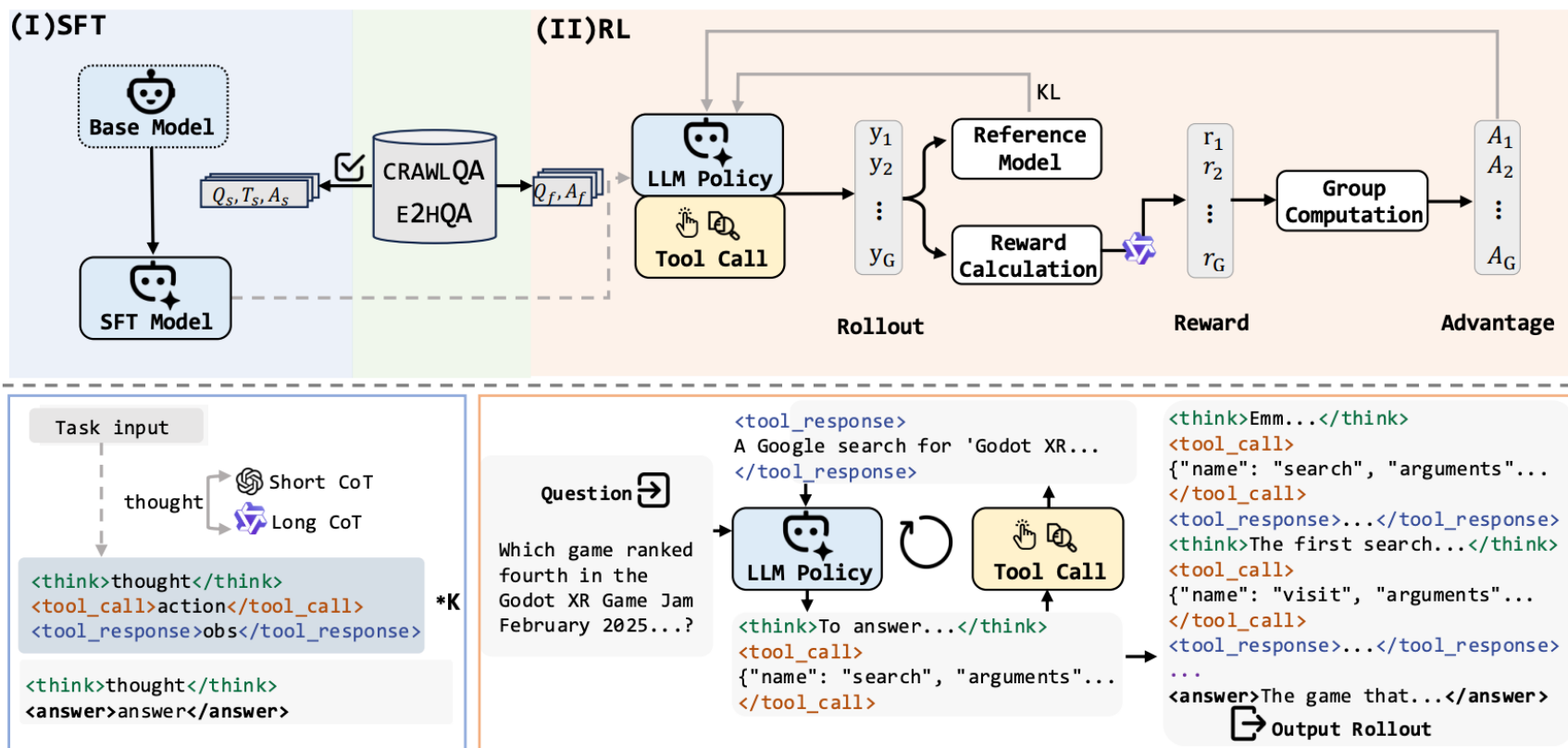
Mimic human behavior by systematically clicking and collecting subpages accessible through sublinks.

Datasets



Rewrite simple questions into more complex, challenging ones systematically .

Methods



Demos

Chatbot

On June 6, 2023, an article by Carolyn Collins Petersen was published in Universe Today. This article mentions a team that produced a paper about their observations, linked at the bottom of the article. Find this paper. Under what NASA award number was the work performed by R. C. Arendt supported by?


点击上传

提交

Agents

请选择一个 Agent

Agent@SEARCH



Agent@SEARCH

我是通用搜索场景下的智能体 (SEARCH)，能够边搜索边思考，欢迎试用!!!

插件

☒ search ☒ visit

推荐对话

列出3个国家及其首都

杭州西湖到杭州西站怎么走

出一份三天两夜的2025年端午北京旅游攻略

对比下最新小米汽车和保时捷性能参数，然后根据最终的结果分析下性价比最高的车型，并给出杭州的供应商

量子计算突破对现有加密体系的威胁

人工智能伦理框架的全球差异

老龄化社会对全球养老金体系的长期冲击

全球碳中和目标下的能源转型路径差异

塑料污染在海洋食物链中的累积效应

AI生成内容 (如AI绘画) 对传统艺术价值的重构

通过 API 使用 · 使用 Gradio 构建 · Settings

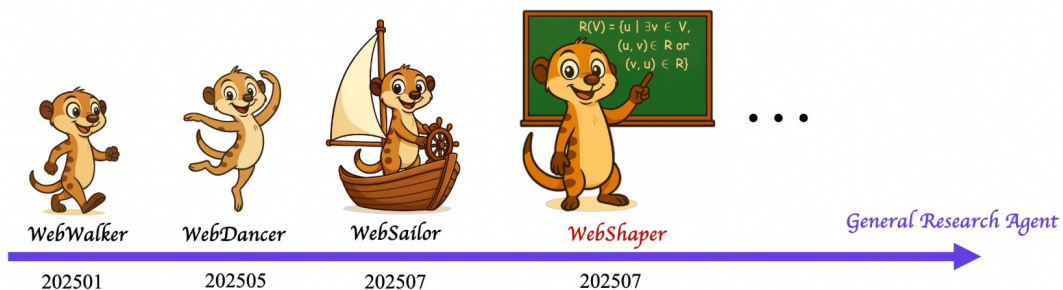
WebAgent

WebAgent for Information Seeking built by Tongyi Lab, Alibaba Group 

1 GITHUB TRENDING
#1 Repository Of The Day

 [WebSailor-3B](#) |  [ModelScope WebSailor-3B](#) |

 [WebDancer-QwQ-32B](#) |  [ModelScope WebDancer-QwQ-32B](#) |  [WebWalkerQA](#)



You can check the paper of [WebDancer](#) and [WebWalker](#) and [WebSailor](#) and [WebShaper](#).

1 GITHUB TRENDING
#1 Repository Of The Day

<https://github.com/Alibaba-NLP/WebAgent>

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

Thanks for watching!
QA