

Theory of Agent: From Definition, to Behavior and Objective

(Toward a Theory of Agents as Tool-Use Decision-Makers)

Dr. Hongru WANG

<https://hrwise-nlp.github.io/>

University of Edinburgh



Content

- ❑ **New Agent Framework Inspired by Cognitive Science: Tool-use Decision Maker**
 - ❑ **Definition, Behavior and Objective of Agent**
- ❑ Theory of Agent Inspires LLM/Agent Principles, Opportunities and Challenges
 - ❑ Principle 1: Foundation --- Self-awareness for Knowledge Boundary and Decision Boundary
 - ❑ Principle 2: Uniqueness and Diversity --- Every LLM/Agent has a Unique Knowledge/Decision Boundary.
 - ❑ Principle 3: Dynamic Conservation --- “Combination” of Intelligence
 - ❑ A Recipe for Agentic Pretraining / SFT / RL / Prompting
- ❑ Future Direction and Conclusions
 - ❑ 1. The Scaling of both Reasoning and Acting
 - ❑ 2. The Scaling of both Agent and Env (or World Model)
 - ❑ 3. The Scaling of Time: Self-evolving Agent
 - ❑ Conclusions

Agents as the Last Mile of Intelligence to Users

Introducing deep research

An agent that uses reasoning to synthesize large amounts of online information and complete multi-step research tasks for you. Available to Pro users today, Plus and Team next.

Try on ChatGPT ↗

Compile a research report on how the retail industry has changed over the past 3 years. Use bullets and tables where appropriate.

Could you specify which aspects of the retail industry you're most interested in?

- E-commerce vs. brick-and-mortar trends
- Consumer behavior shifts
- Supply chain challenges
- Emerging technologies (AI, automation, etc.)

Manus AI

Home

Manus AI Cases

Request Invitation Code

OpenAI Deep F

Manus AI - The AI Assistant That Turns Your Thoughts Into Actions

Manus AI is a world-leading general-purpose AI agent designed to help users efficiently complete various complex tasks. The name Manus comes from the Latin word for 'hand,' symbolizing its ability to execute tasks. It has achieved state-of-the-art (SOTA) performance across all three difficulty levels in the GAIA benchmark, far surpassing other AI assistants.

Get Started with Manus AI →

Request Invitation Code

Manus

Manus Joins Meta for Next Era of Innovation



OSWorld: Benchmarking Multimodal Agents for Open-Ended Tasks in Real Computer Environments

Yao Fu¹, Danyang Zhang¹, Jixuan Chen¹, Xiaochuan Li¹,
Yuhao Zhou¹, Jiahao Jiang¹, Zhoujun Cheng¹, Dongchan Shin¹, Fangyu Lei¹, Yitao Liu¹,
Yue Hou³, Silvio Savarese², Caiming Xiong², Victor Zhong⁴, Tao Yu¹
¹Salesforce Research, ²Carnegie Mellon University, ⁴University of Waterloo

OpenAI Data Data Viewer Slides Twitter Discord

Computer-Using Agent

GAIA Leaderboard

With augmented capabilities due to added tooling, efficient prompting, access to search, etc). (See our [FAQ](#) for more details.)

GAIA requires different levels of tooling and autonomy to solve. It is therefore divided in 3 levels, where level 1 should be breakable by very good LLMs, and level 3 indicate a strong jump in model capabilities. In private answers and metadata.

1. JSON. Some questions come with an additional file, that can be found in the same folder and whose id is given in the field `file_name`.

2. Models.

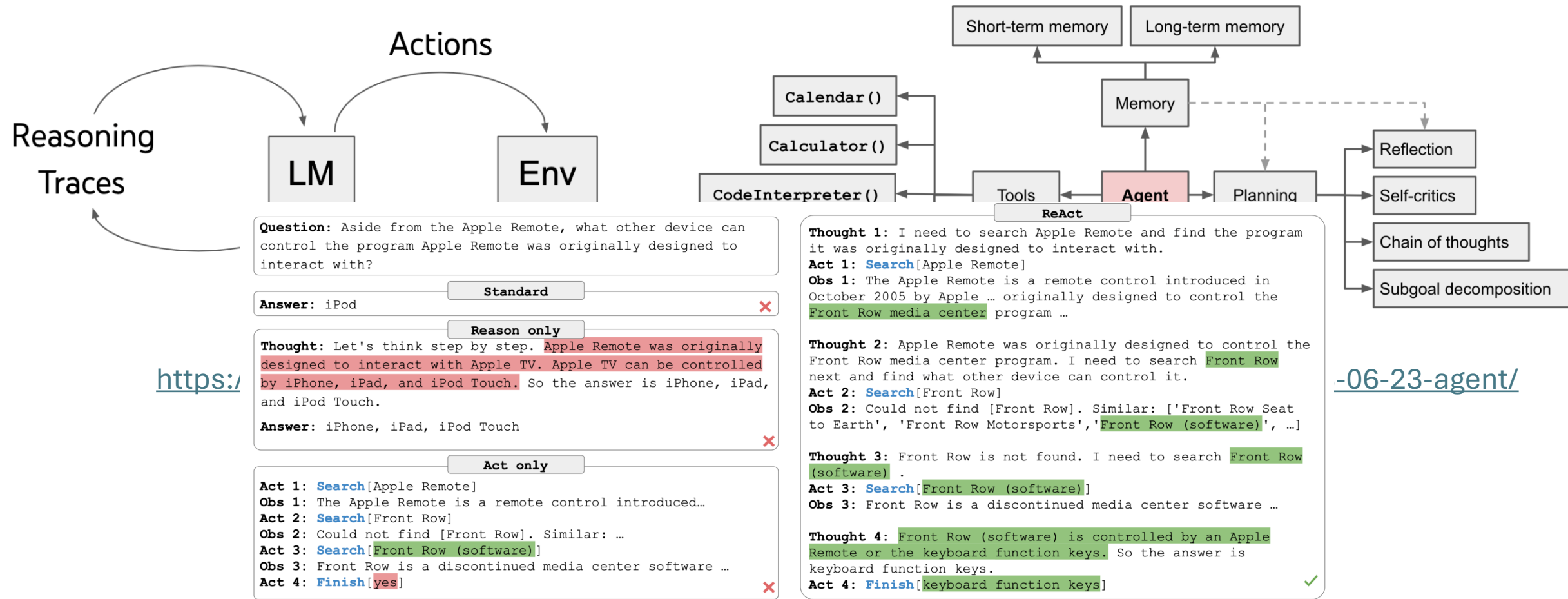
Average scores over different runs when possible in our paper, we only report the best run in the leaderboard.

	organisation	Average score (%)	Level 1 score (%)	Level 2 score (%)	Level 3 score (%)
		99.39	98.11	100	100
		99.39	98.11	100	100
	Princeton AI Lab	87.27	88.68	89.53	79.92
		87.27	96.23	90.7	57.69
	Princeton AI Lab	86.06	96.23	86.05	65.38
		83.83	92.45	87.21	50
	Skywork AI	82.42	92.45	83.72	57.69
	Skywork AI	88	92.45	79.07	57.69
		79.39	88.68	80.23	57.69
	Princeton AI Lab	78.79	88.68	79.07	57.69
		78.18	86.79	77.91	65.54
		77.58	90.57	75.58	57.69

Refresh

Alita reaches top 1 at GAIA (validation, 2025.6)

Previous / Existing Popular Agent Definitions



<https://>

-06-23-agent/

Agent = [Reasoning + Acting] * n

The Relationship Between Reasoning and Acting

Reasoning and acting are

- Different tokens for model
- Different tools / actions for the agent
- Different interactions for user
-

But they are epistemic equal means to acquire knowledge to solve the task.



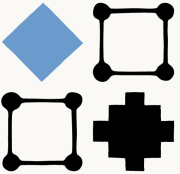
Shunyu Yao
@ShunyuYao12

To reason and act is the same thing

翻译帖子

下午10:56 · 24/6/24 来自 Earth · 865 次查看

Engineering at Anthropic



The "think" tool: Enabling Claude to stop and think in complex tool use situations

Published Mar 20, 2025 A new tool that improves Claude's complex problem-solving performance

```
{
  "name": "think",
  "description": "Use the tool to think about something. It will not obtain new information or change the database, but just append the thought to the log. Use it when complex reasoning or some cache memory is needed.",
  "input_schema": {
    "type": "object",
    "properties": {
      "thought": {
        "type": "string",
        "description": "A thought to think about."
      }
    },
    "required": ["thought"]
  }
}
```

Google Scholar

cognitive tool

Articles About 5,860,000 results (0.13 sec)

Any time
Since 2025
Since 2024
Since 2021
Custom range...

Sort by relevance
Sort by date

Any type
Review articles

☐ include patents
☒ include citations

Create alert

What are cognitive tools?

DH Jonassen - **Cognitive tools** for learning, 1992 - Springer
... **tools** that extend the mind This workshop was about **cognitive tools** - computer-based **tools** ... Computer-based **cognitive tools** are in effect **cognitive** amplification **tools** that are part of ...
☆ Save 📄 Cite Cited by 508 Related articles All 5 versions

[PDF] Technology as cognitive tools: Learners as designers

DH Jonassen - ITForum Paper, 1994 - tecfa.unige.ch
... **Cognitive tools** are generalizable computer **tools** that ... **Cognitive tools** and environments activate **cognitive** learning strategies and critical thinking. They are computationally based **tools** ...
☆ Save 📄 Cite Cited by 383 Related articles All 4 versions 🔗

[book] Computers as Cognitive Tools: 1

SP Lajoie, SJ Derry - 1993 - books.google.com
... are employed, and the forms of "**cognitive tools**" that are embedded within systems to help ... computers as **tools** for enhancing learning. Computers as **Cognitive Tools** is appropriate for ...
☆ Save 📄 Cite Cited by 924 Related articles All 10 versions 🔗

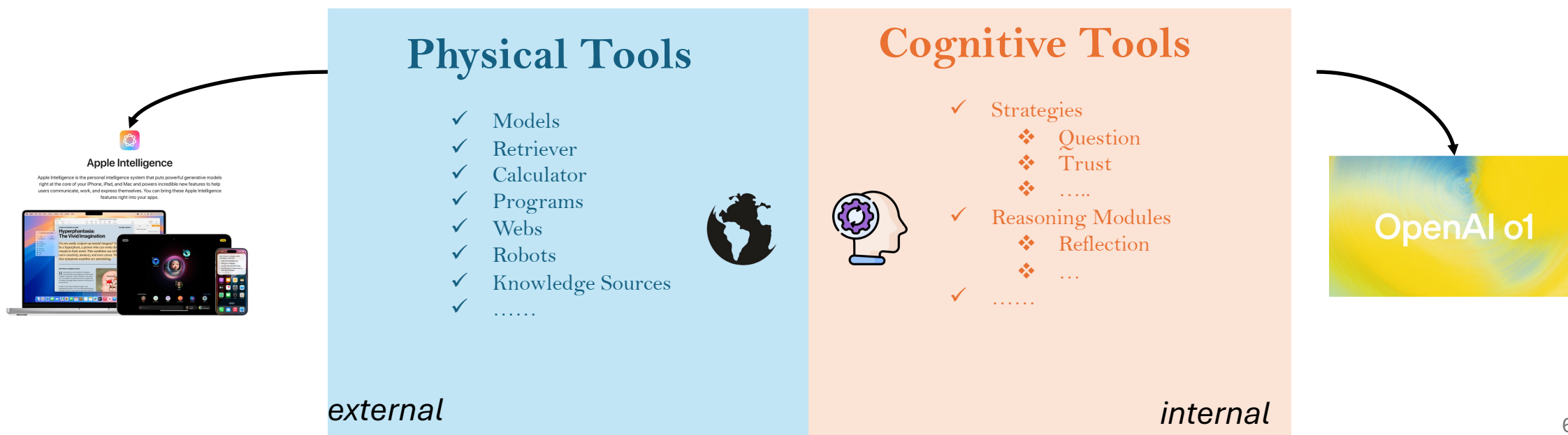
[book] Cognitive tools for learning

PAM Kommers, DH Jonassen, JT Mayes - 1992 - research.utwente.nl
... to address the theme of **cognitive tools** as discussed in this book ... **tools** and was the main reason that '**cognitive tools**' became ... during instruction allows for **cognitive** amplification. Some ...
☆ Save 📄 Cite Cited by 342 Related articles All 8 versions 🔗

<https://ysymyth.github.io/The-Second-Half/>
<https://www.anthropic.com/engineering/claude-think-tool>

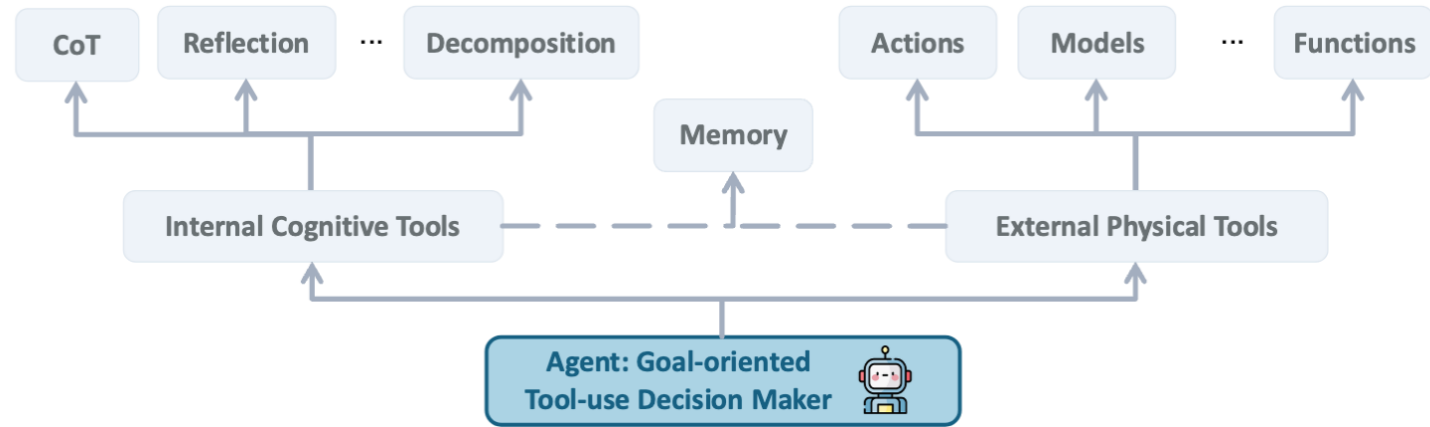
Reasoning \sim Acting = Tools

- Tool is defined as object that can extend an individual's ability to modify features of the surrounding environment or help them accomplish a particular task in general. It can be **internal cognitive/conceptual tools** (i.e., *reasoning*) and **external physical tools** (i.e., *acting*).
 - Internal cognitive/conceptual tool** refer to specifies an internal cognitive mechanisms that aids systematic or investigative thought, to retrieve internal knowledge of agent about current state (e.g, **internal world model**).
 - External physical tool** refer to external modules that are invoked by a rule or a specific token and whose outputs are incorporated into the context of agent (e.g., **external world model**).

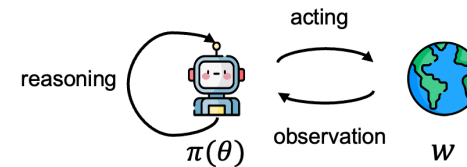


New Agent Definition

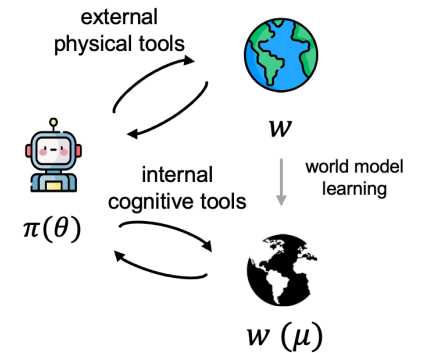
- An agent is an entity that **coordinates** internal cognitive tools (e.g., reflection) and external physical tools (e.g., function callings) to acquire knowledge in order to achieve a specific goal.



- Internal cognitive tools and external physical tools are **epistemic equal means** to acquire knowledge to solve the task, as shown in Figure (b).

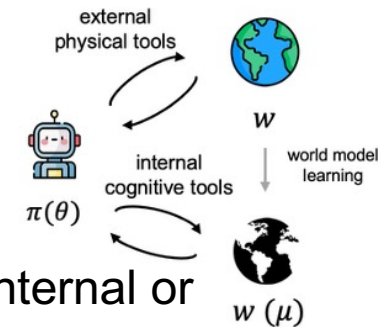


(a) ReAct-based Agent



(b) Tool-integrated Agent

Three Advantages of Tool-integrated Agents



- **Unified Format:** $\tau = (t_1, k_1, t_2, k_2, \dots, t_n, k_n)$

- t_n, k_n stands for tool call and returned knowledge at n_{th} step. The tool could be either internal or external.

- **Flexible and Robust**

- It degrades to previous ReAct paradigm if we consider the internal tools and the whole reasoning part, then it becomes $(r_1, t_1, k_1, \dots, r_n, t_n, k_n)$ here t_n, k_n only

- If we solely consider internal tools, it is proved that simply outcome-based reward tool utilization such as reflection and decomposition to solve the problem in LLMs (i.e., DeepSeek-R1). Alternatively, simply outcome-based reward also triggers tool utilization as evidenced in recent studies (i.e., Search-R1, ToRL, OTC-PO).



Percy Liang
@percyliang

What is the analogue of next-token prediction for reinforcement learning? To get true generality, you want to be able to convert everything in the world to an environment+reward for training.

翻译帖子

下午10:50 · 27/2/25 · 5.4万 次查看

22

33

293

180



- **Potential Next Scaling Law**

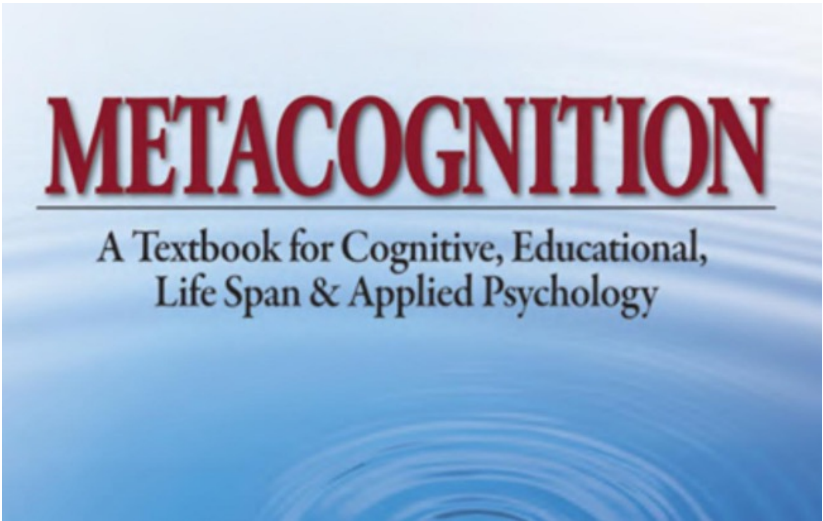
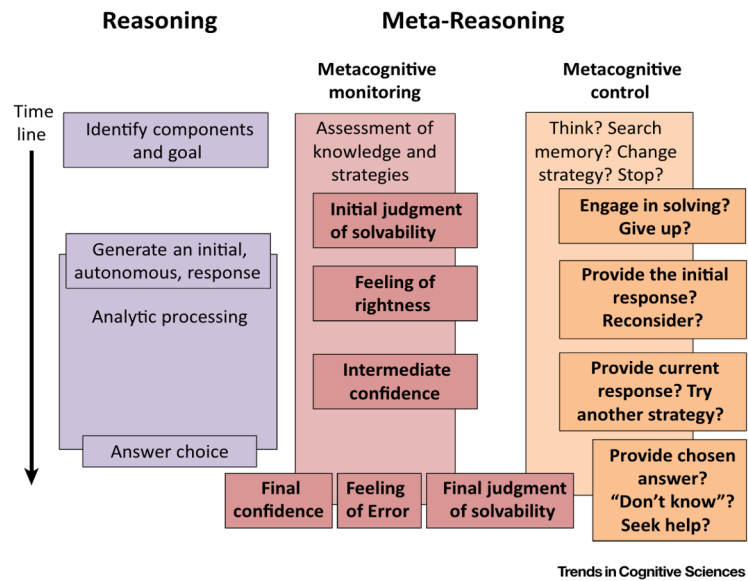
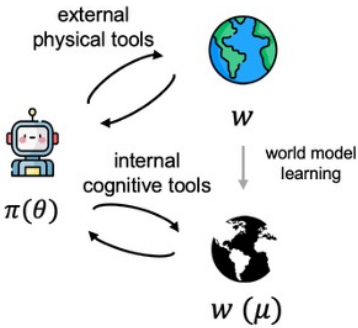
- **Next Tool/State Prediction:** Just as next-token prediction enables LLMs to learn a compressed representation of the world from text, next-tool prediction allows agents to learn procedural knowledge through interaction.
 - Procedural knowledge scaling \propto Context / experience scaling, leading to self-evolving agent

Content

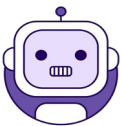
- ❑ **New** Agent Framework Inspired by Cognitive Science: Tool-use Decision Maker
 - ❑ Definition, Behavior and Objective of Agent
- ❑ **Theory of Agent Inspires LLM/Agent Principles, Opportunities and Challenges**
 - ❑ Principle 1: Foundation --- Self-awareness for Knowledge Boundary and Decision Boundary
 - ❑ Principle 2: Uniqueness and Diversity --- Every LLM/Agent has a Unique Knowledge/Decision Boundary.
 - ❑ Principle 3: Dynamic Conservation --- “Combination” of Intelligence
 - ❑ A Recipe for Agentic Pretraining / SFT / RL / Prompting
- ❑ Future Direction and Conclusions
 - ❑ 1. The Scaling of both Reasoning and Acting
 - ❑ 2. The Scaling of both Agent and Env (or World Model)
 - ❑ 3. The Scaling of Time: Self-evolving Agent
 - ❑ Conclusions

How to **coordinates** internal and external tools?

- How **human** call different tools in mind: meta-reasoning theory, metacognition,



- How **agent** call different tools?
 - The key also lies in monitoring and control



Monitoring ?

Judgement of solvability
Intermediate confidence
Reward model
Uncertainty estimation

...

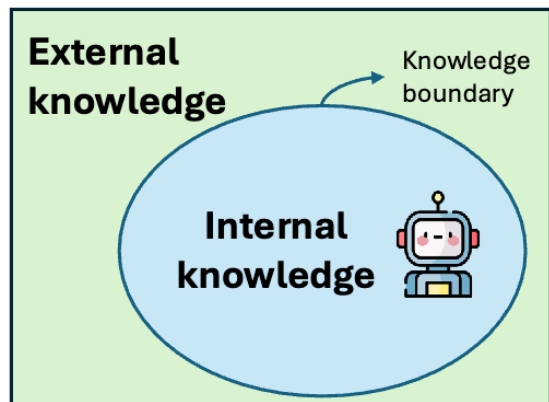
Control ?

Cognitive tools
Physical tools
...

Meta-Reasoning: Monitoring and Control of Thinking and Reasoning

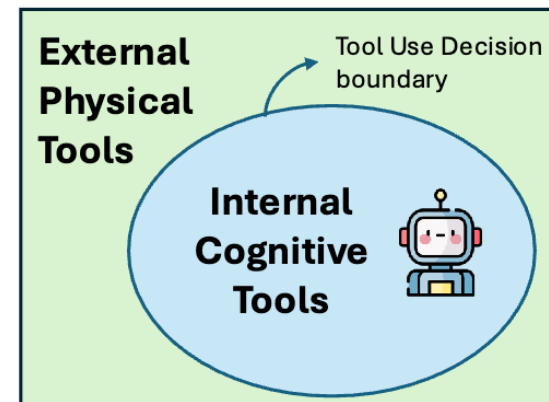
Two Concepts for Monitoring and Control

- To make thing easier, let's assume all knowledge is correct, and can be accessed via tools, and there is a way to accurately identify the knowledge boundary.



Monitor: Self-aware Knowledge Boundary

Decides
→



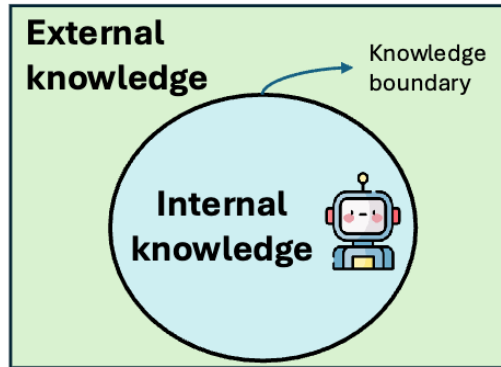
Control: Self-aware Tool Utilization

- We hope that LLMs can **utilize internal cognitive tools to gain internal knowledge** while **only call external tools to gain external knowledge** during problem-solving processing. (explain later)
 - The challenge here is **self-aware tool utilization**

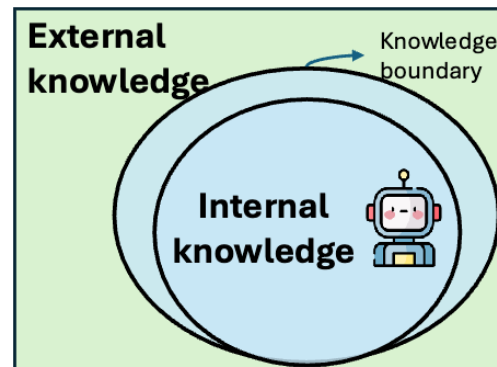
Optimize Tool Use Decision Boundary to match Knowledge Boundary (知行合一)

Principle 1: Foundations

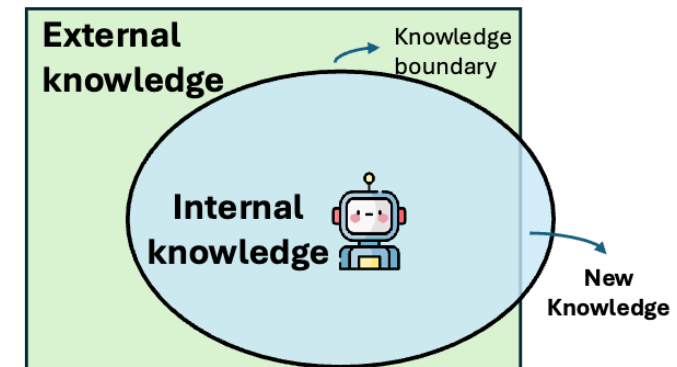
- **Lemma 1.1:** *Generally*, as time advances, the model's capabilities evolve and the knowledge boundary expands.
- **Lemma 1.2:** *Specifically*, the knowledge/decision boundaries can be redistributed, e.g., through continual training, allowing for strengthening in specific domains.



(a) Knowledge Boundary



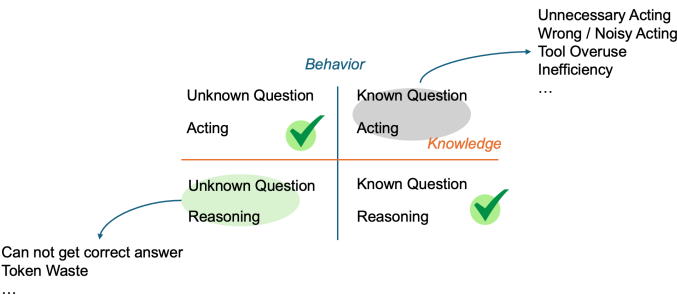
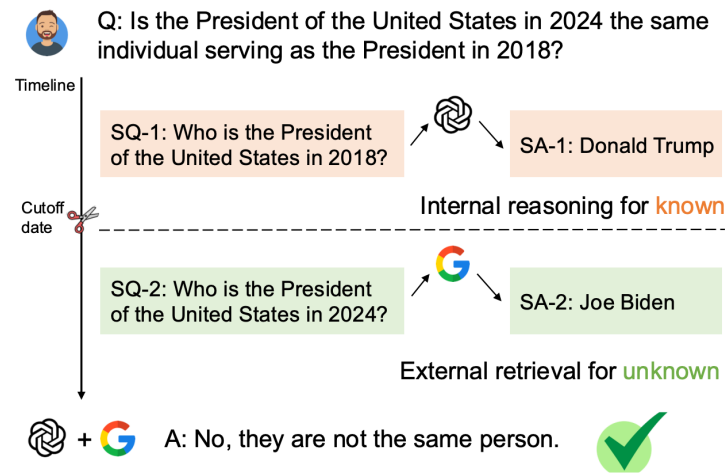
(b) Knowledge Expansion



(c) New Knowledge Discovery

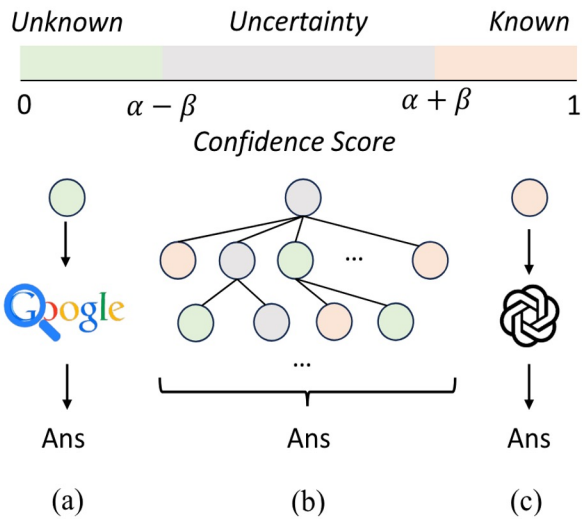
Self-DC: When to Reason and When to Act?

- **Single Known.** The question contains no sub-questions and can be solved using internal knowledge of LLMs, such as with the generate-then-read method.
- **Single Unknown.** The question contains no sub-questions and can only be solved using external knowledge, such as with the retrieve-then-read method.
- **Compositional Known.** The question contains several sub-questions, and each sub-question is *Single Known*.
- **Compositional Unknown.** The question contains several sub-questions, and at least one sub-question is *Single Unknown*.



Better alignment between two boundaries brings better trade-off between effectiveness and efficiency.

First Compositional unknown Question Answering dataset (CuQA)



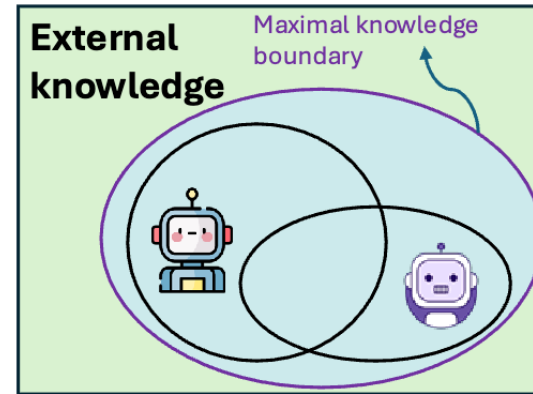
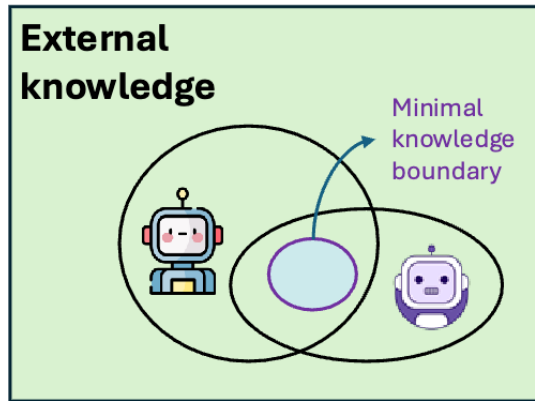
- Defining reasoning and acting as different functions / tools. *Call these functions leveraging model-based planning, and meta-reasoning theory (confidence scores).*
- Solving compositional/complex problems in *different level of granularity.*
- **Simple and Scalable** Purely based on self-aware capabilities of LLMs. As LLM evolves, the framework evolves.

Methods	#R	CuQA			FreshQA		
		EM	F1	Acc [†]	EM	F1	Acc [†]
<i>w/o retrieval</i>							
Direct	0	29.0	19.4	46.4	27.2	17.3	53.0
CoT	0	28.8	18.2	46.0	29.2	18.1	53.8
Few-shot-CoT*	0	43.0	3.2	50.8	35.0	9.1	55.4
GenRead	0	29.6	29.2	47.4	26.8	27.7	52.0
<i>w/ retrieval</i>							
RR	n	32.0	31.6	55.4	35.2	32.6	63.4
REFEED	$2n$	26.2	<u>33.5</u>	51.8	28.8	<u>34.5</u>	57.4
IRCoT	$3n$	47.8	13.5	64.6	34.2	17.8	61.4
Self-Ask*	$0-n$	19.8	3.8	48.4	5.6	9.8	59.0
ITER-RETGEN*	$2n$	23.4	12.6	50.9	31.2	21.1	55.8
<i>Self-DC (verb)</i>	$0-n$	34.0	32.2	53.8	30.2	30.2	59.8
<i>Self-DC (prob)</i>	$0-n$	<u>36.4</u>	36.5	<u>56.4</u>	37.4	36.6	66.4

First Framework to consider relationship between reasoning and acting

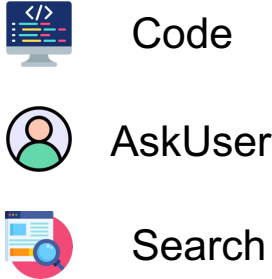
Principle 2: Uniqueness and Diversity

- **Lemma 2.1:** Each model has its own knowledge boundary and decision boundary.
- **Lemma 2.2:** There exist minimal and maximal knowledge (and decision) boundaries across *all* models.

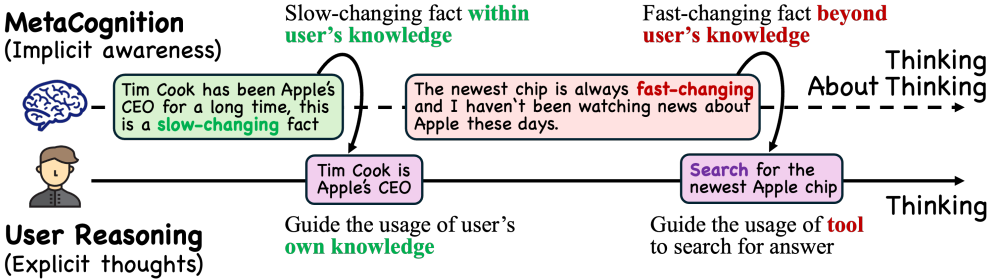


SMART: Self-Aware Agent for Tool Overuse Mitigation

- We adapt three established dataset to create the meta-reasoning chain:
 - Math: *simple arithmetic* v.s. **challenging calculation** (e.g., MATH)
 - Intention: *commonsense* v.s. **user specific intentions** (e.g., Intention-in-Interaction)
 - Time: *never-changing facts* v.s. **fast-changing facts** (e.g., FreshQA)



What is newest chip developed by the company whose CEO is Tim Cook?



SMARTAgent achieves **higher accuracy with lower tool call number and higher confidence in decision**

One-fit-for-all strategy is approximating **Maximal Knowledge Boundary (lemma 2.2)**

Method	Model	Math (MATH)		Time (FreshQA)		Intention (Intention-in-Interaction)		
		Tool Used ¹ (Times)	Accuracy ² (%)	Tool Used ¹ (Times)	Accuracy ² (%)	Tool Used ¹ (Times)	Missing Details Recovery ³ (Lv3 / Lv2, %)	Summarized Intention Coverage ⁴ (%)
Open-Source								
Normal Reasoning Trained	Mistral-7B	0.00	17.00	0.00	48.00	0.00	41.86 / 43.84	-
	Llama-3.1-8B	0.00	41.00	0.00	48.00	0.00	38.37 / 42.49	-
Base Model Reasoning Prompt	Mistral-7B	0.00	17.25	0.00	29.00	0.00	37.21 / 33.06	-
	Llama-3.1-8B	0.00	53.00	0.00	26.00	0.00	40.70 / 25.76	-
	Mistral-Nemo(12B)	0.00	47.00	0.00	33.00	0.00	44.19 / 28.37	-
	Mistral-Small(24B)	0.00	72.25	0.00	34.00	0.00	41.86 / 31.82	-
	Llama-3.1-70B	0.00	70.00	0.00	36.00	0.00	41.86 / 29.24	-
Base Model Tool Prompt	Mistral-7B	3.90	13.25	1.67	49.00	3.80	48.84 / 21.70	63.04
	Llama-3.1-8B	1.93	51.00	2.05	56.00	3.77	54.76 / 25.90	70.20
	Mistral-Nemo(12B)	2.35	46.00	1.19	59.00	1.80	31.35 / 5.82	59.27
	Mistral-Small(24B)	1.55	76.00	1.73	62.00	2.52	45.74 / 33.62	78.20
	Llama-3.1-70B	3.53	67.50	2.78	63.00	2.71	45.74 / 35.96	61.68
SMARTAgent	Mistral-7B	0.60 _(3.30)	22.75 _(15.50)	1.00 _(6.67)	64.00 _(15.00)	3.60 _(10.20)	74.42 _(25.58) / 65.44 _(21.60)	81.76 _(18.72)
	Llama-3.1-8B	0.88 _(11.05)	54.75 _(11.75)	1.05 _(11.00)	67.00 _(11.00)	3.80 _(10.03)	81.40 _(26.61) / 67.41 _(24.92)	78.28 _(18.08)
	Mistral-Nemo(12B)	0.82 _(11.53)	49.50 _(12.50)	1.00 _(10.19)	70.00 _(11.00)	3.34 _(11.54)	77.91 _(33.72) / 62.15 _(23.78)	82.30 _(23.03)
	Mistral-Small(24B)	0.79 _(10.76)	69.75 _(16.25)	1.00 _(10.73)	66.00 _(11.00)	3.89 _(11.37)	74.42 _(38.68) / 68.87 _(35.25)	84.99 _(26.79)
	Llama-3.1-70B	0.94 _(12.59)	72.50 _(12.50)	1.01 _(11.07)	66.00 _(11.00)	3.51 _(10.80)	68.60 _(22.86) / 58.15 _(22.19)	86.09 _(24.41)
	Tool Used Macro-Average Decrease (%)		24.00		Performance Macro-Average Increase (%)		37.10	
Closed-Source								
Base Model Reasoning Prompt	GPT-4o-mini	0.00	73.00	0.00	44.00	0.00	45.35 / 32.41	-
	GPT-4o	0.00	79.50	0.00	47.00	0.00	38.37 / 28.54	-
Base Model Tool Prompt	GPT-4o-mini	2.55	54.50	1.06	56.00	1.91	50.00 / 26.90	76.44
	GPT-4o	0.27	79.25	1.01	65.00	1.17	40.70 / 15.61	86.80

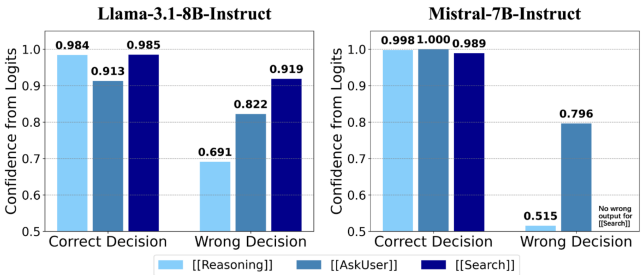


Figure 5: Confidence analysis shows that SMART effectively enhances the model’s decision-making confidence in selecting the correct reasoning approaches.

Principle 3: Dynamic Conservation of Knowledge

- **Lemma 3.1:** At any time step t , the total world knowledge W_t is fixed and identical across all models.
- **Lemma 3.2:** For any task or query q and model m , there exists a minimal and fixed epistemic effort $N(q, m)$ allocated between internal and external sources, that is necessary to solve the task, such as $N(q, m) = K_{int} + K_{ext}$.
 - **Task-Model dependency Optimization:** $N(q, m)$ is jointly determined by the complexity of the task and the capabilities of the model.
 - **Capability Equivalence via Dynamic Offloading:** Even models with limited internal capacity can achieve same performance by dynamically offloading reasoning or retrieval steps to more capable tools or agents. There is no difference between 8B ($K_{ext} \rightarrow N$) and 70B ($K_{int} \rightarrow N$) from Agent perspective considering models as one of tools.
 - **Agent Objective:** Pursuing the optimal behavior that minimize interactions while managing latency, cost, and constraints, besides the final correctness.

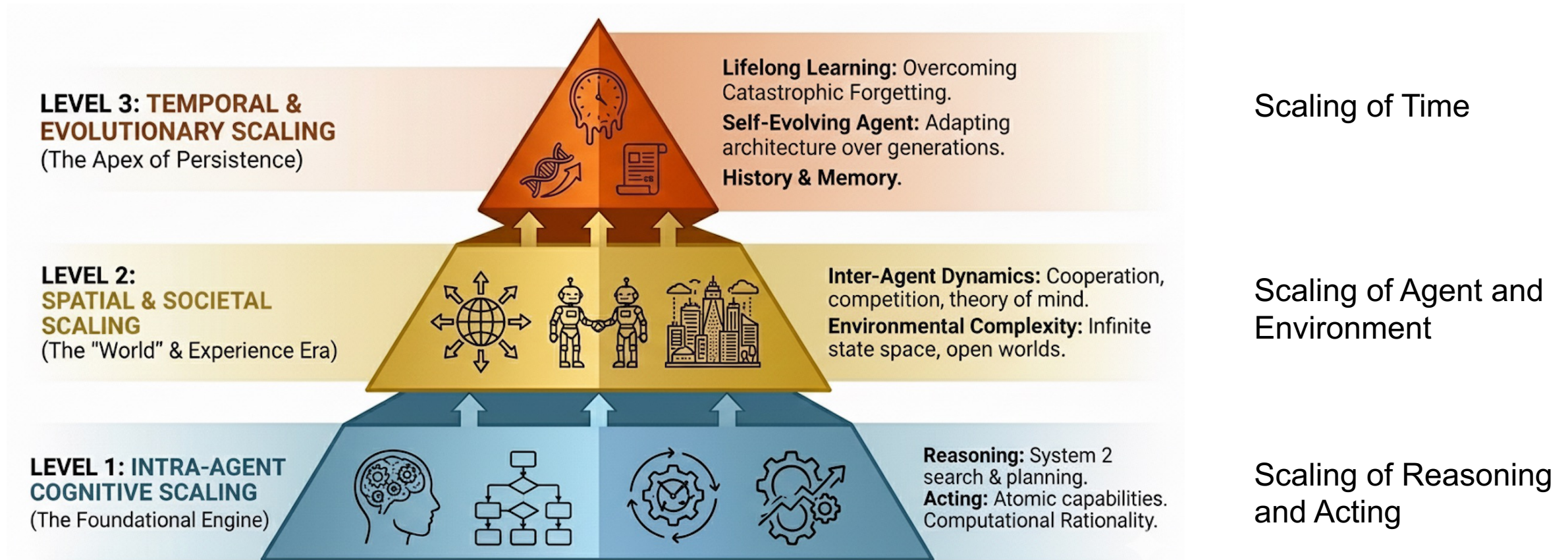
A Roadmap to Autonomous Agent

- **Agentic Pretraining:** Next tool prediction, As research trends toward unified agent architectures, modeling all forms of interaction (API calls, UI navigation, or environment manipulation) as structured, learnable outputs opens the door to a new kind of scaling law: one that governs knowledge acquisition, not just compression.
 - Unified Format: $\tau = (t_1, k_1, t_2, k_2, \dots, t_n, k_n)$
 - Data Collection: It is extremely challenging to collect massive pretraining interaction corpus (only Big Companies)
- **Agentic Supervised-finetuning:** It is important to collect model-task-specific trajectories instead of collecting one trajectory for all models due to lemma 2.1. Additionally, it is more effective to leverage the lemma 2.2 by utilizing maximal knowledge boundary to build one-fits-all dataset.
- **Agent Reinforcement Learning:** Reinforcement learning (RL) offers a more promising path for aligning a model's decision-making with its own knowledge boundary, as agents can learn from experience how to adaptively use tools. The key challenge lies in designing reward functions that go beyond correctness.
- **Agent Prompting:** Once the model is trained, previous numerous studies utilize prompt engineering to develop task-specific agentic workflows across various domains. Despite achieving exceptional performance on complex tasks, few of these approaches rigorously evaluate behavioral optimality, such as internal cognitive tool overuse (i.e., overthinking) or external physical tool overuse (i.e., overacting).

Content

- ❑ **New** Agent Framework Inspired by Cognitive Science: Tool-use Decision Maker
 - ❑ Definition, Behavior and Objective of Agent
- ❑ Theory of Agent Inspires LLM/Agent Principles, Opportunities and Challenges
 - ❑ Principle 1: Foundation --- Self-awareness for Knowledge Boundary and Decision Boundary
 - ❑ Principle 2: Uniqueness and Diversity --- Every LLM/Agent has a Unique Knowledge/Decision Boundary.
 - ❑ Principle 3: Dynamic Conservation --- “Combination” of Intelligence
 - ❑ A Recipe for Agentic Pretraining / SFT / RL / Prompting
- ❑ **Future Direction and Conclusions**
 - ❑ 1. The Scaling of both Reasoning and Acting
 - ❑ 2. The Scaling of both Agent and Env (or World Model)
 - ❑ 3. The Scaling of Time: Self-evolving Agent
 - ❑ Conclusions

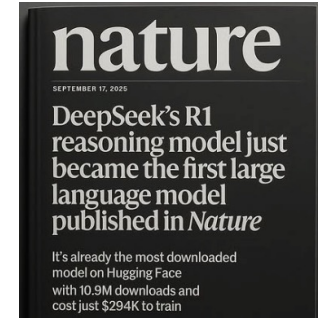
Three Levels of Scaling Simultaneously



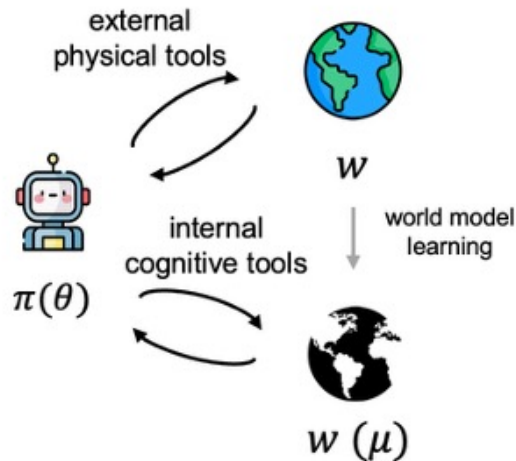
Future Direction 1: Scaling of both Reasoning and Acting

RL has been proven effective in scaling reasoning (i.e., DeepSeek-R1) and acting (i.e., Kimi-K2) capabilities, respectively.

However, how to scale them together without losing any part of capabilities?



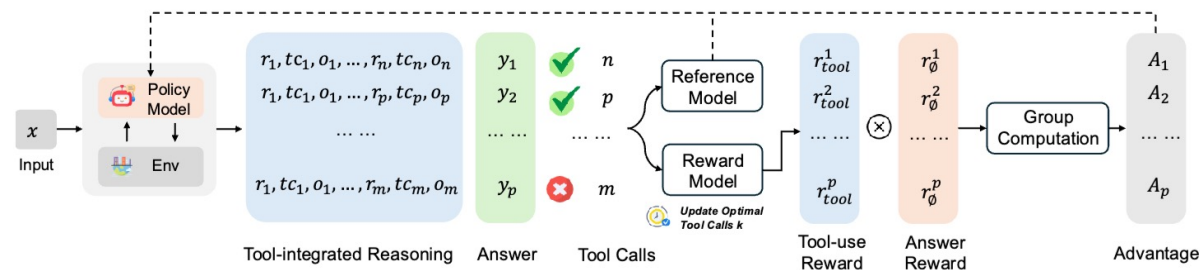
*"The autonomous machine intelligence is designed to **minimize the number of actions** a system needs to take in the real world to learn a task. It does so by learning a world model that capture as much knowledge about the world as possible without taking actions in the world."* --- Yann Lecun [1]



This is both the goal itself and a means of achieving it.

As long as the agent can complete the task successfully, minimizing external physical tools means maximizing the internal tools by our **Principle 3**, also means the agent can internalize the external environment partly or fully (i.e., learn a better internal world model).

Acting Less is Reasoning More (OTC-PO)

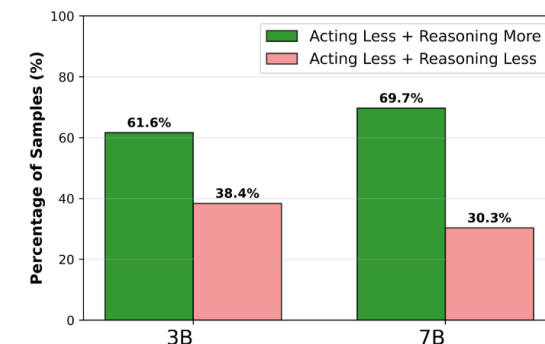
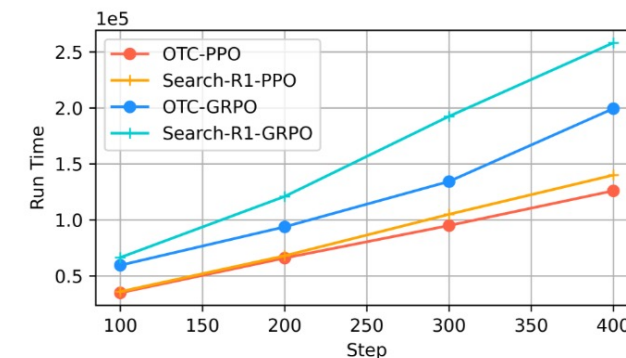


$$r_{\phi}^{tool}(q, y) = \alpha * r_{tool} * r_{\phi}(q, y)$$

A simple, faster, and generalizable **OTC-PO** algorithm to encourage the model to use fewer tool calls to solve the problem

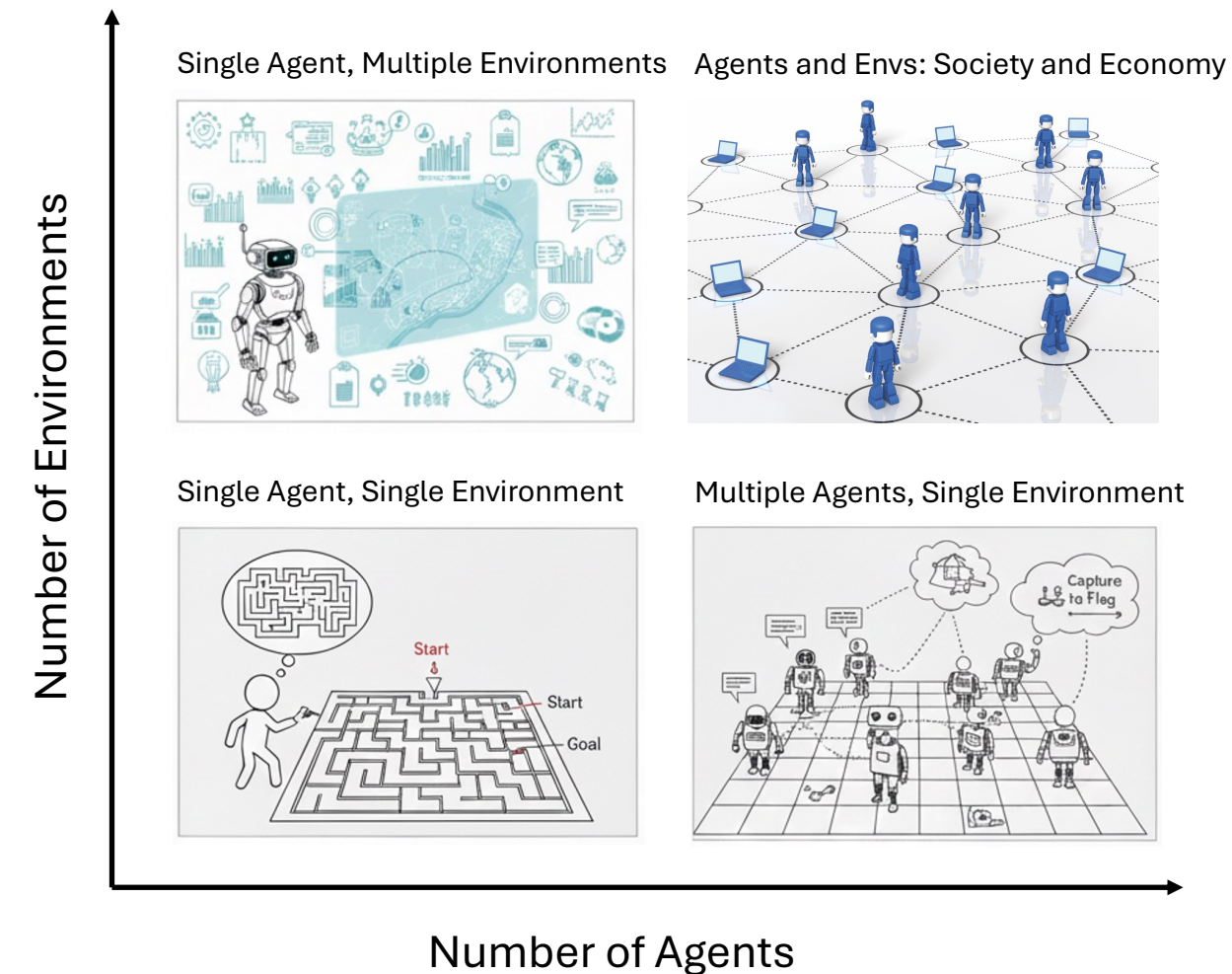
Question: Are Giuseppe Verdi and Ambroise Thomas both Opera composers ? Answer: Yes	
Search-R1-PPO: <think> I need to determine if Giuseppe Verdi and Ambroise Thomas are both opera composers. </think> <search> Giuseppe Verdi </search> <observation> ... </observation> <search> Ambroise Thomas </search> <observation> ... </observation> <search> Are Giuseppe Verdi and Ambroise Thomas both Opera composers? </search> <observation> ... </observation> <think> Based on the observations, both Giuseppe Verdi and Ambroise Thomas are described as composers, specifically opera composers. </think> <answer> yes </answer> (Tool Calls = 3)	Search-R1-GRPO: <think> To answer this question, I need to check if both Giuseppe Verdi and Ambroise Thomas are composers of opera. I don't have immediate access to this knowledge, but I can conduct a search to find out. </think> <search> who are Giuseppe Verdi and Ambroise Thomas and their works </search> <observation>...</observation> <think>From the information gathered, ... This confirms Thomas as an opera composer. Next, I will search on Giuseppe Verdi. </think> <search> who is Giuseppe Verdi </search> <observation>...</observation> <think>Giuseppe Verdi is named as an Italian opera composer of prominent stature. This confirms Verdi as an opera composer. Answering the initial question: Question: Are Giuseppe Verdi and Ambroise Thomas both Opera composers? </think> <answer> yes </answer> (Tool Calls = 2)
OTC-PPO: <think> First, I need to determine if Giuseppe Verdi and Ambroise Thomas are both opera composers. </think> <search> Giuseppe Verdi and Ambroise Thomas both Opera composers </search> <observation> ... </observation> <think> Based on the search results, Giuseppe Verdi and Ambroise Thomas are both known as opera composers. </think> <answer> yes </answer> (Tool Calls = 1)	OTC-GRPO: To answer the question of whether Giuseppe Verdi and Ambroise Thomas are both Opera composers, let's first reason through the information: 1. Giuseppe Verdi is a well-known Italian opera composer. He is famous for his numerous operas, including "Rigoletto," "La Traviata," and "Aida." 2. Ambroise Thomas is a French opera composer. He is known for his operas such as "Mignon" and "Hamlet." Based on this reasoning, both Giuseppe Verdi and Ambroise Thomas are indeed Opera composers. <answer> yes </answer> (Tool Calls = 0)

Models	NQ			HotpotQA		
	EM (↑)	TC (↓)	TP (↑)	EM (↑)	TC (↓)	TP (↑)
Qwen2.5-3B(-Base)						
R1-Base	0.226	-	-	0.201	-	-
SFT	0.249	-	-	0.186	-	-
RAG	0.348	1.0	0.348	0.255	1.0	0.255
IRCoT	0.111	10.0	0.011	0.164	10.0	0.016
Search-R1-PPO	0.403	1.738	0.232	0.279	1.716	0.163
Search-R1-GRPO	0.404	1.426	0.283	0.312	1.802	0.173
OTC-PPO	0.355	1.010 (▼ 41.9%)	0.351 (▲ 51.3%)	0.260	1.026 (▼ 40.2%)	0.253 (▲ 55.2%)
OTC-GRPO	0.444	1.008 (▼ 29.3%)	0.440 (▲ 55.5%)	0.365	1.387 (▼ 23.0%)	0.263 (▲ 52.0%)
Qwen2.5-7B(-Base)						
R1-Base	0.270	-	-	0.242	-	-
SFT	0.318	-	-	0.217	-	-
RAG	0.349	1.0	0.349	0.299	1.0	0.299
IRCoT	0.224	9.999	0.022	0.133	9.982	0.013
Search-R1-PPO	0.449	3.282	0.136	0.380	3.741	0.102
Search-R1-GRPO	0.399	1.697	0.235	0.341	2.109	0.162
OTC-PPO	0.446	1.040 (▼ 68.3%)	0.429 (▲ 215.4%)	0.383	1.464 (▼ 60.9%)	0.262 (▲ 156.9%)
OTC-GRPO	0.444	0.990 (▼ 41.7%)	0.448 (▲ 90.6%)	0.366	1.005 (▼ 52.3%)	0.364 (▲ 124.7%)



*Less tool calls,
less time,
less money, but
more reasoning,
more intelligence,
more scalable.*

Future Direction 2: Scaling of both Agent and Env



Scaling Environments for LLM Agents in the Era of Learning from Interaction: A Survey

Yuchen Huang[♡] Sijia Li[♡] Minghao Liu[♡] Wei Liu[♣] Zhiyuan Fan[♡]
Yi R. (May) Fung[♡]
[♡]Hong Kong University of Science and Technology
[♣]King's College London
{yhuanggn, yrfung}@cse.ust.hk

Google Research Google DeepMind

Towards a Science of Scaling Agent Systems

Yubin Kim^{1,3,†}, Ken Gu¹, Chanwoo Park³, Chunjong Park², Samuel Schmidgall², A. Ali Heydari¹, Yao Yan¹, Zhihan Zhang¹, Yuchen Zhuang², Yun Liu¹, Mark Malhotra¹, Paul Pu Liang³, Hae Won Park³, Yuzhe Yang¹, Xuhai Xu¹, Yilun Du¹, Shwetak Patel¹, Tim Althoff¹, Daniel McDuff¹ and Xin Liu^{1,†}

¹Google Research, ²Google DeepMind, ³Massachusetts Institute of Technology, [†]Corresponding Author

Most of existing work scale either agents or environments. However, **how to scale both of them together remains under-explored.**

Word2World: Can LLM be implicit Text-based World Models?

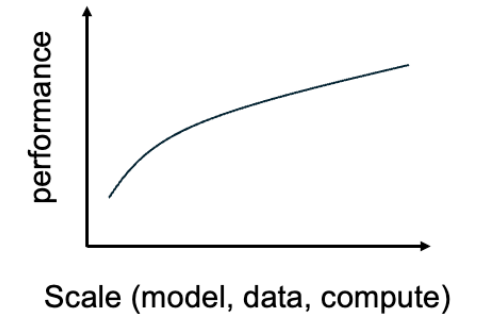
Agent (LLM)

Next token predication
In-context Learning
Scaling law

....

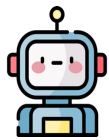
Env (World Model)

Master 'worlds' by learning structured,
predictive representations of
environments.



Can we identify a similar path to guide the development of world model?

Text-based Environments as Bridge
(i.e., AlfWorld, SciWorld, WebShop, ...)



Agent (LLM)

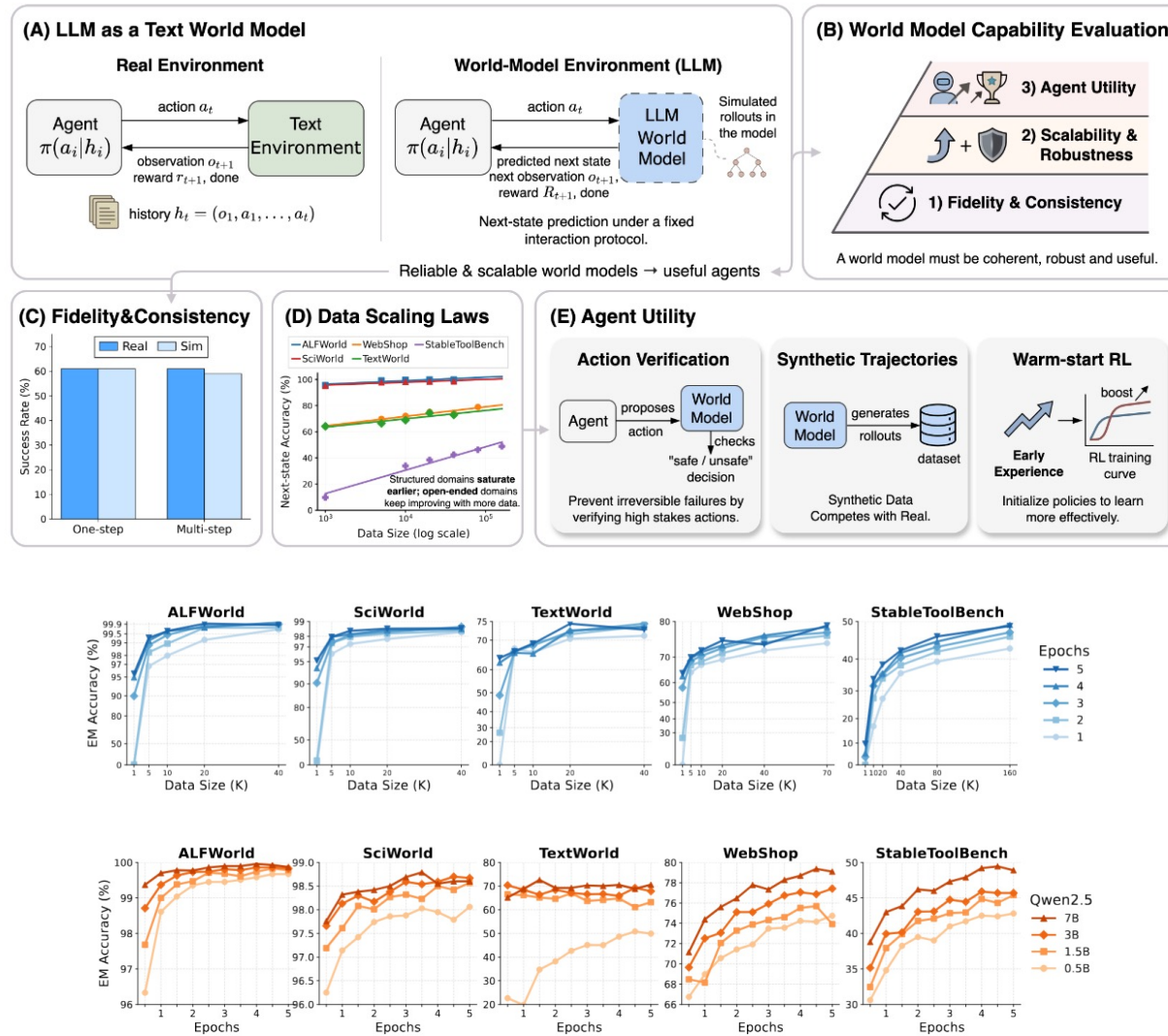


Env (World Model)

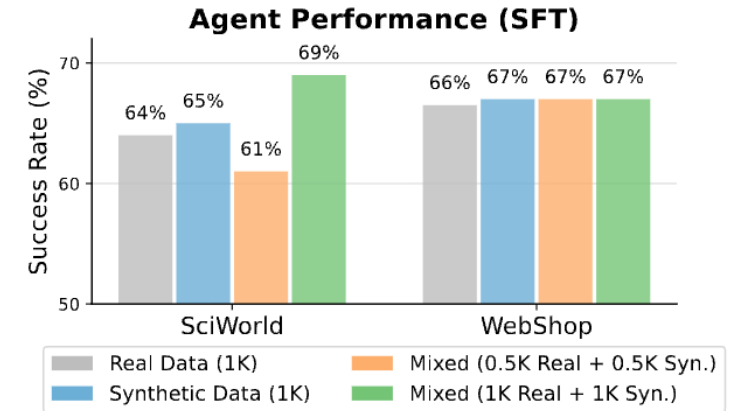
Next states prediction $(s, a) \rightarrow s'$
In-context Learning $(s, a, s', a') \rightarrow s''$
World model scaling law

....

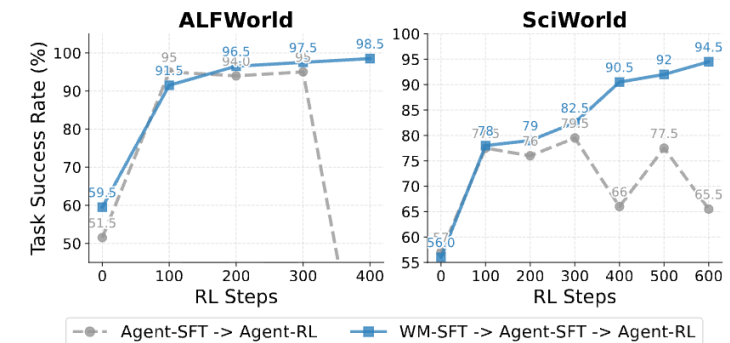
Word2World: Can LLM be implicit Text-based World Models?



Yes, simple fine-tuning unlocks near-perfect state prediction and long-horizon consistency



Data Synthesis



World model performance scales predictably with data and model size, mirroring LLMs. However, the nature of this scaling is tied to environment complexity.

Warm-start RL

Future Direction 3: Scaling of Time

Welcome to the Era of Experience

David Silver, Richard S. Sutton*

Abstract

We stand on the threshold of a new era in artificial intelligence that promises to achieve an unprecedented level of ability. A new generation of agents will acquire superhuman capabilities by learning predominantly from experience. This note explores the key characteristics that will define this upcoming era.

Published in Transactions on Machine Learning Research (01/2026)

A Survey of Self-Evolving Agents

What, When, How, and Where to Evolve on the Path to Artificial Super Intelligence

Huan-ang Gao[†], Jiayi Geng[†], Wenye Hua[†], Mengkang Hu[†], Xinzhe Juan[†], Hongzhang Liu[†], Shilong Liu[†], Jiahao Qiu[†], Xuan Qi[†], Qihan Ren[†], Yiran Wu[†], Hongru Wang[†], Han Xiao[†], Yuhang Zhou[†], Shaokun Zhang[†], Jiayi Zhang[†], Jinyu Xiang, Yixiong Fang[†], Qiwen Zhao[†], Dongrui Liu[†], Cheng Qian[†], Zhenhailong Wang[†], Minda Hu[†], Huazheng Wang[†], Qingyun Wu[†], Heng Ji[†], Mengdi Wang[†]

[†]Princeton University, [‡]Princeton AI Lab, [§]Tsinghua University, [¶]Carnegie Mellon University, [⋄]University of Sydney, [⋆]Shanghai Jiao Tong University, [⋈]Pennsylvania State University, [⋊]University of Michigan, [⋋]Oregon State University, [⋌]The Chinese University of Hong Kong, [⋍]Fudan University, [⋎]The Hong Kong University of Science and Technology (Guangzhou), [⋏]The University of Hong Kong, [⋐]University of California, Santa Barbara, [⋑]University of California San Diego, [⋒]University of Edinburgh, [⋓]University of Illinois Urbana-Champaign

Github Repo: <https://github.com/CharlesQ9/Self-Evolving-Agents>

[†]Equal contribution and the order is determined alphabetically, [⋈]Corresponding Author

Reviewed on OpenReview: <https://openreview.net/forum?id=CTr3bovS5F>

First comprehensive survey about self-evolving agent

LifelongAgent-2026

CALL FOR PAPERS

First Workshop on Lifelong Agent: Learning-Aligning-Evolving

An ICLR 2026 Workshop

2026
26th/Apr
Full day
Rio de Janeiro

Building Agents That Endure: Adaptive Learning, Stable Alignment, Sustainable Growth.

01 Speakers & Panellists

Siva Reddy
McGill / Mila

Azalia Mirhoseini
Stanford / Deepmind

Su Yu
OSU / NeoCognition

Graham Neubig
CMU / All Hands AI

Asli Celikyilmaz
Meta

Sergey Levine
UC Berkeley

Manos Koukoumidis
Oumi AI

02 Organizers & Advisors

Cheng Qian
UIUC

Emre Can Acikgoz
UIUC

Hongru Wang
University of Edinburgh

Zhenfei Yin
Oxford

Manling Li
Northwestern

Vivian Chen
National Taiwan University

Guanhua Chen
SUSTech

Jiahao Qiu
Princeton

Caiming Xiong
Salesforce AI Research

Heng Ji
UIUC

Gokhan Tur
UIUC

Dilek Hakkani-Tür
UIUC

Philip Torr
Oxford

Kam-Fai Wong
CUHK

Jun Wang
UCL

Mengdi Wang
Princeton

First workshop focus on lifelong agent: learning, aligning and evolving (call for paper)

<https://lifelongagent.github.io/>

Alita, Alita-G, Agent-Distill,

ALITA: GENERALIST AGENT ENABLING SCALABLE AGENTIC REASONING WITH MINIMAL PREDEFINITION AND MAXIMAL SELF-EVOLUTION

Jiahao Qiu^{*1}, Xuan Qi^{*2}, Tongcheng Zhang^{*3}, Xinzhe Juan^{3,4}, Jiacheng Guo¹, Yifu Lu¹, Yimin Wang^{3,4}, Zixin Yao¹, Qihan Ren³, Xun Jiang², Xing Zhou³, Dongrui Liu³, Ling Yang², Yue Wu⁴, Kaixuan Huang³, Shilong Liu¹, Hongru Wang⁵, Mengdi Wang¹

¹AI Lab, Princeton University ²IIS, Tsinghua University ³Shanghai Jiao Tong University
⁴University of Michigan ⁵Tianqiao and Chrissy Chen Institute ⁶The Chinese University of Hong Kong

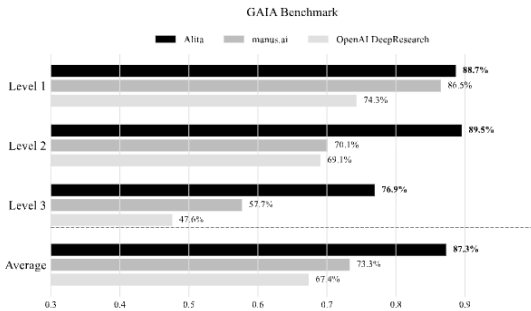


Figure 1: Performance of Alita, manus.ai, and OpenAI DeepResearch[1]

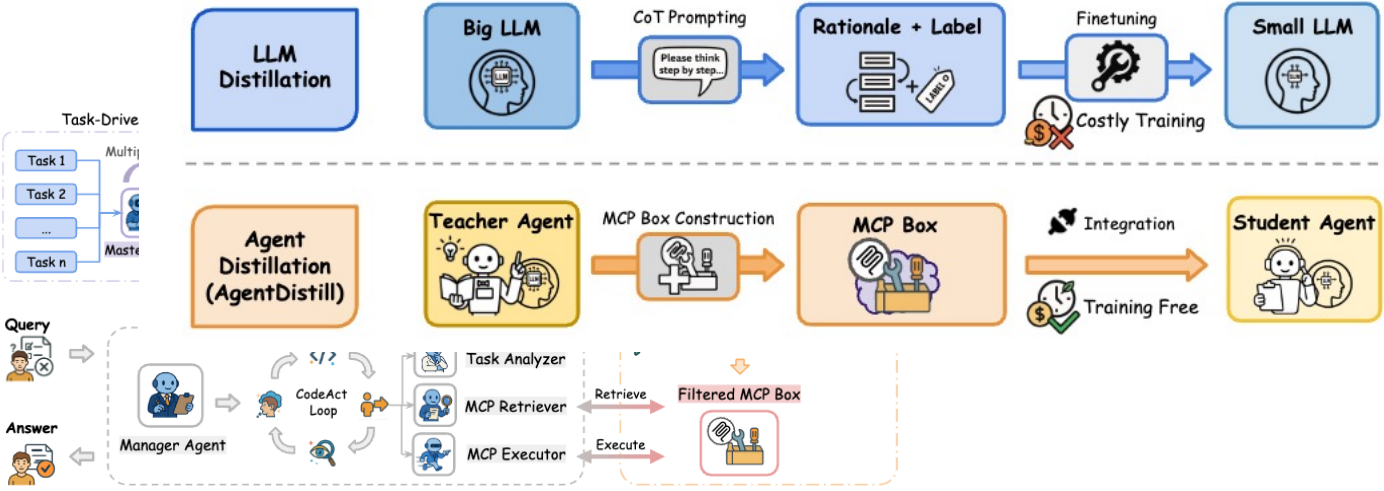


Alita can create

AGENTDISTILL: TRAINING-FREE AGENT DISTILLATION WITH GENERALIZABLE MCP BOXES

Jiahao Qiu^{*1}, Xinzhe Juan^{*2,4}, Yimin Wang^{*2,4}, Ling Yang^{*1}, Xuan Qi³, Tongcheng Zhang⁴, Jiacheng Guo¹, Yifu Lu¹, Zixin Yao⁵, Hongru Wang⁶, Shilong Liu¹, Xun Jiang⁷, Liu Leqi^{†8}, Mengdi Wang^{†1}

¹AI Lab, Princeton University ²University of Michigan ³IIS, Tsinghua University
⁴Shanghai Jiao Tong University ⁵Columbia University ⁶The Chinese University of Hong Kong
⁷Tianqiao and Chrissy Chen Institute ⁸University of Texas at Austin



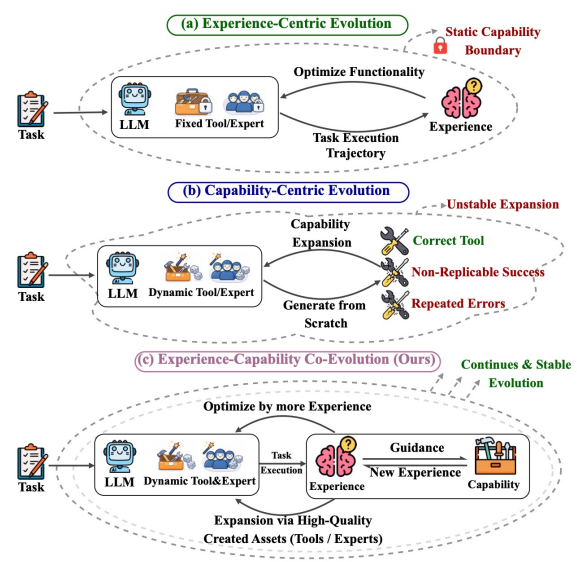
ALITA-G: SELF-EVOLVING GENERATIVE AGENT FOR AGENT GENERATION

Jiahao Qiu^{*1}, Xuan Qi^{*2}, Hongru Wang^{*1,3}, Xinzhe Juan^{4,5}, Yimin Wang^{4,5}, Zelin Zhao⁶, Jiayi Geng¹, Jiacheng Guo¹, Peihang Li⁷, Jingzhe Shi², Shilong Liu^{1,8}, Mengdi Wang^{1,8}

¹Princeton University ²Tsinghua University ³The Chinese University of Hong Kong
⁴Shanghai Jiao Tong University ⁵University of Michigan ⁶King's College London ⁷Hong Kong University

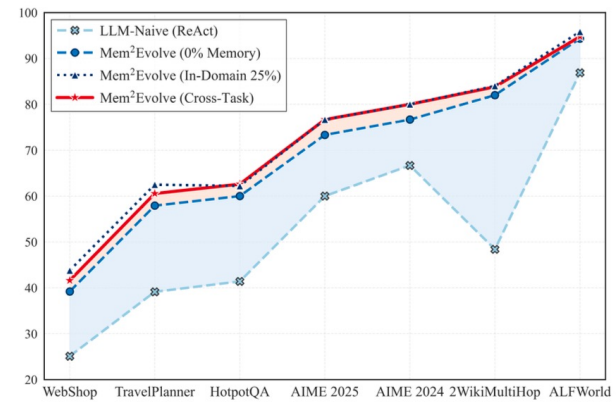
Alita-G can create specialized agent automatically during evolution.

Mem²Evolve: Co-evolution of Agents and Envs



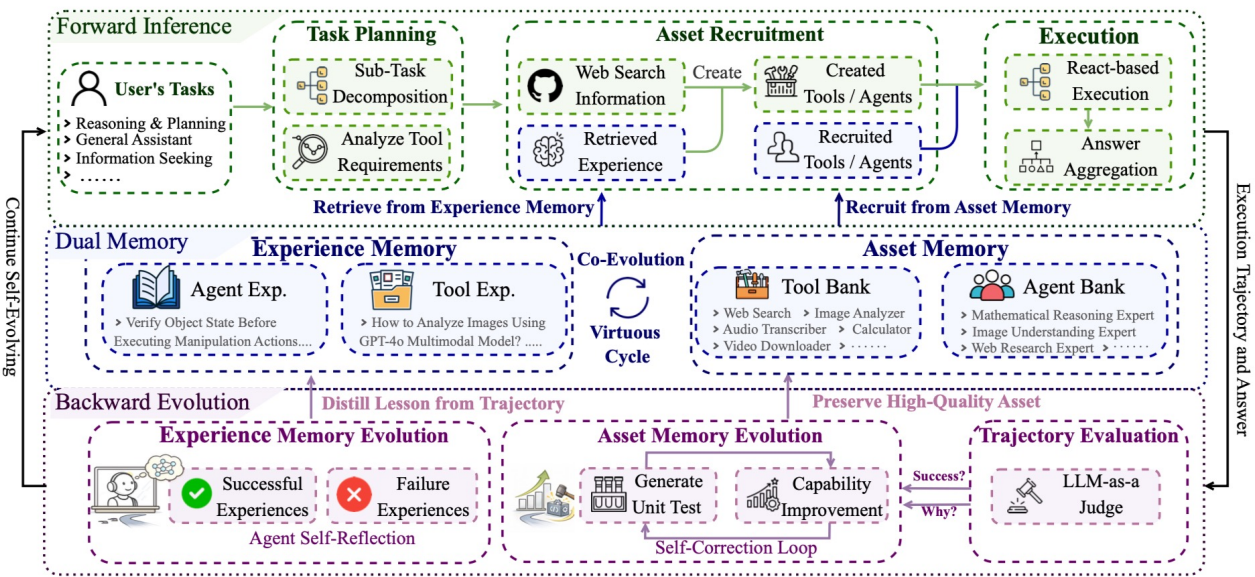
Framework	Experience Distillation			Capability Expansion				Exp.-Guided Creation
	Optimization	Persistence	Source	Tool Crea.	Agent Crea.	Tool/Agent	Crea. Grounding	
DSpy (Khattab et al., 2023)	✓	✗	🧠	✗	✗	Static	—	✗
DyLAN (Liu et al., 2023)	✗	✓	🧠	✗	✗	Static	—	✗
ReasoningBank (Ouyang et al., 2025)	✗	✓	🧠	✗	✗	Static	—	✗
AFlow (Zhang et al., 2025a)	✓	✗	🧠	✗	✗	Static	—	✗
AgentSquare (Shang et al., 2025)	✓	✗	🧠	✗	✗	Static	—	✗
Agentic Neural Networks (Ma et al., 2025)	✓	✗	🧠	✗	✗	Static	—	✗
AgentVerse (Chen et al., 2023)	✓	✗	—	✗	✓	Dynamic	🧠	✗
AutoAgents (Chen et al., 2024)	✓	✗	—	✗	✓	Dynamic	🧠	✗
SwarmAgentic (Zhang et al., 2025b)	✓	✗	—	✗	✓	Dynamic	🧠	✗
Alita (Qiu et al., 2025)	✓	✗	—	✓	✗	Dynamic	🧠 + 🌐	✗
ToolMaker (Wölflin et al., 2025)	✗	✗	—	✓	✗	Dynamic	🧠 + 🌐	✗
Mem ² Evolve (Ours)	✓	✓	🧠 + 🌐	✓	✓	Dynamic	🧠 + 🌐 + 🧠	✓

Table 1: Comparison of self-evolving agent frameworks. **Optimization** indicates whether experience is used to optimize the agent (e.g., prompts). **Persistence** denotes whether experiences are persistently stored for future reuse. **Source**: 🧠 agent task execution trajectory, 🌐 tool creation process. **Tool Crea.** and **Agent Crea.** indicate whether the framework supports creation of tools and expert agents, respectively. **Tool/Agent** denotes whether the toolset and expert agents are static or dynamic. **Crea. Grounding** indicates the knowledge sources used for asset creation, 🧠 parametric knowledge, 🌐 web search information, 🧠 experience. **Exp.-Guided Creation** indicates whether new assets are created under the guidance of past experience. Details in the Appendix A.1 and A.2.



Cross-task evolution

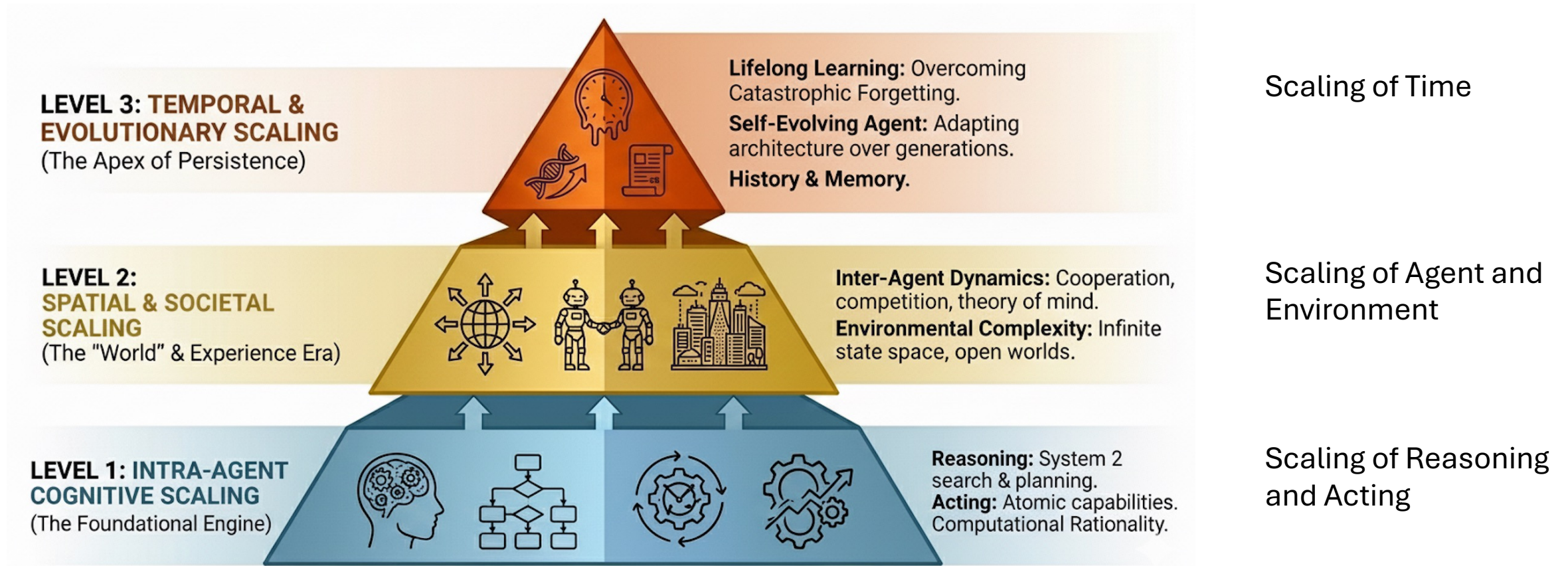
Three ways of evolutions



Method	GAIA				Embodied	Multi-Hop QA			Math		Planning	Web Interaction		Avg.
	L1	L2	L3	Total	ALFWorld	HotpotQA	2Wiki	AIME24	AIME25	TravelPlanner		WebShop		
Naive-Large Language Model														
GPT-5-Chat (Direct)	16.98	12.79	7.69	12.49	83.58	50.40	<u>81.80</u>	60.00	46.67	38.68		22.31		49.49
GPT-5-Chat (CoT)	24.53	17.44	11.54	17.84	83.58	47.40	74.40	66.67	56.67	39.51		27.49		51.71
GPT-5-Chat (ReAct)	26.42	17.44	11.54	18.47	86.87	41.40	48.40	66.67	60.00	39.13		25.10		48.27
OpenAI-DeepResearch [†]	74.29	69.06	<u>47.60</u>	67.36	—	—	—	—	—	—		—		—
Experience-Centric Evolving														
DyLAN	24.53	19.78	11.54	18.62	91.20	52.00	65.00	46.67	43.33	43.15		36.40		49.55
EvoAgent	22.64	19.78	11.54	17.99	92.50	54.40	75.00	66.67	43.33	49.20		37.80		54.61
AFLow	26.42	17.44	15.38	19.75	<u>93.40</u>	60.80	72.40	66.67	63.33	53.24		<u>37.90</u>		58.44
DSPy	30.19	15.12	11.54	18.95	92.80	55.60	76.40	66.67	50.00	44.90		35.50		55.10
Capability-Centric Evolving														
Alita	<u>81.13</u>	<u>75.58</u>	46.15	<u>72.73</u>	86.13	<u>58.80</u>	77.40	<u>70.00</u>	<u>66.67</u>	48.32		30.21		<u>63.78</u>
AgentVerse	30.19	16.28	19.23	21.90	88.32	38.60	74.60	60.00	50.00	47.25		32.53		51.65
AutoAgents	35.85	24.42	19.23	26.50	87.92	54.20	73.80	40.00	36.67	43.52		31.40		49.25
SwarmAgentic	28.30	18.60	13.46	20.40	88.79	56.00	80.00	46.67	40.00	<u>59.14</u>		34.12		53.14
Ours														
Mem ² Evolve	88.68	82.56	57.69	76.31	94.31	60.80	82.00	76.70	73.33	59.25		39.20		70.24

Mem²Evolve Framework

Rethinking *Three* Levels of Scaling Simultaneously



Lots of problem need to be defined and explored, welcome to join theory of agent!

Conclusion

- Agents are not merely an engineering problem; they are becoming increasingly scientific and theoretical, like scaling law of LLMs. We also need to find more theories of agent.
- Agent can be regarded as human. Lots of problem in human society also happens in agent society, i.e., internet/tool addiction.
- Every company should have an agent department instead of LLM department.
- Join theory of agent no matter what you do now, you are not just a researcher, but may next scientist, entrepreneur, and even billionaire.
- Safety and personalization still matters in practice. Aligning decision boundaries with both preference boundaries and knowledge boundaries is tricky.