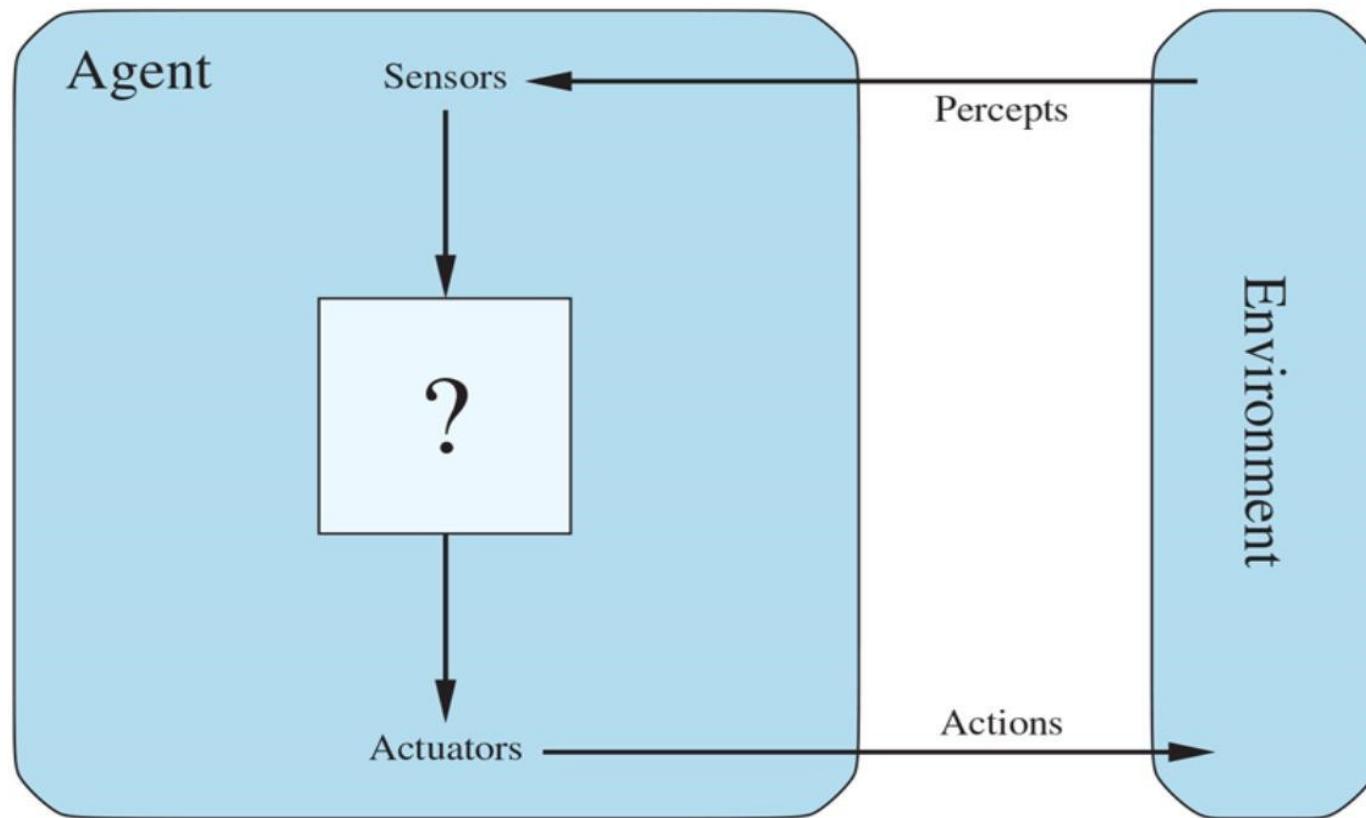


Natural Language Processing

LLM Agents
Instructor: Yangqiu Song

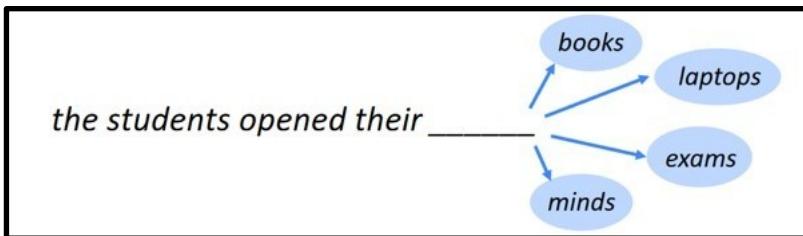
Why agents again?



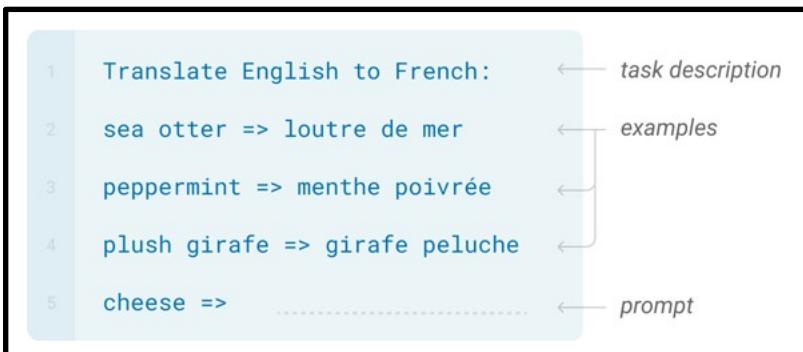
“An **agent** is anything that can be viewed as perceiving its **environment** through **sensors** and acting upon that environment through **actuators**”

-- Russel & Norvig, *AI: A Modern Approach*

The promise of LLMs: Generality and few-shot learning

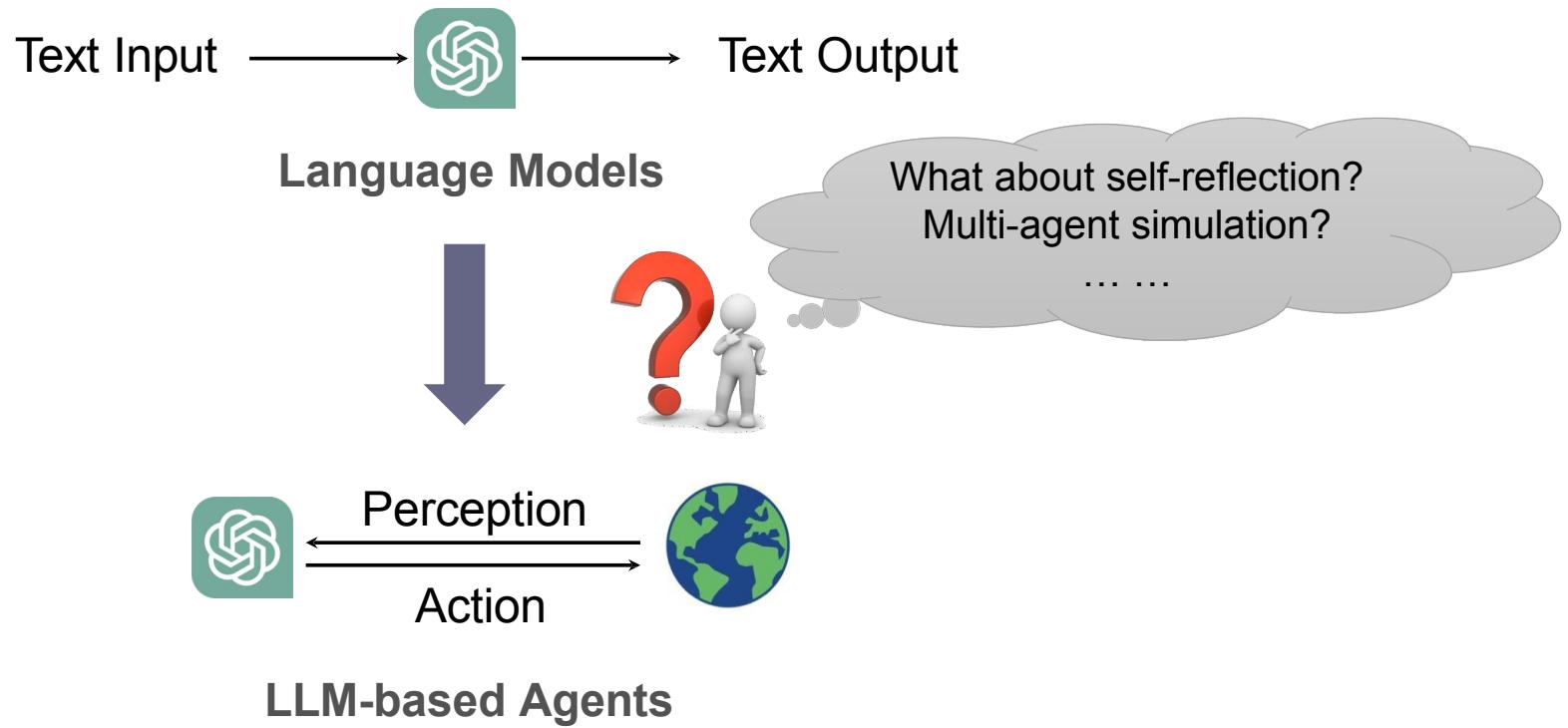


Training: next-token prediction
on massive text corpora



Inference: (few-shot)
prompting for various tasks!

'Modern' agent = LLM + external environment?



Two competing views

LLM-first view We make an LLM into an agent!

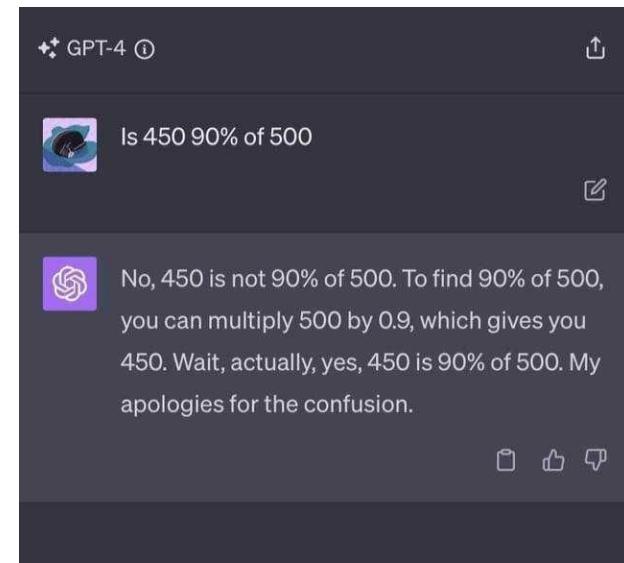
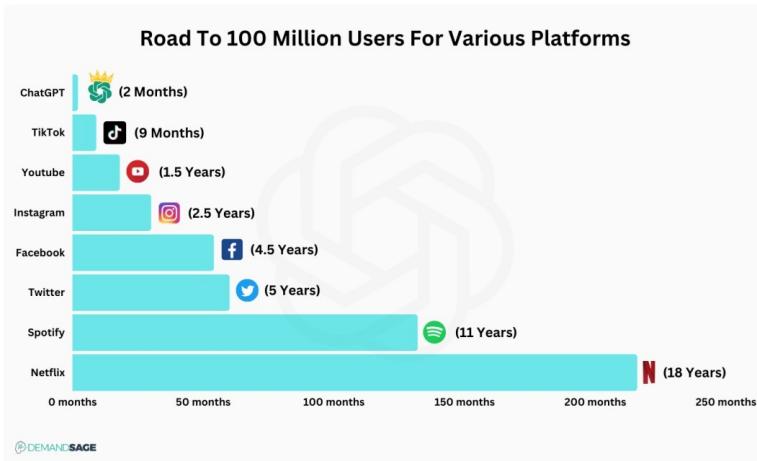
- *Implications:* scaffold on top of LLMs, prompting-focused, heavy on engineering

Agent-first view We integrate LLMs into AI agents so they can use language for reasoning and communication!

- *Implications:* All the same challenges faced by previous AI agents (e.g., perception, reasoning, world models, planning) still remain, but we need to **re-examine them through the new lens of LLMs** and tackle new ones (e.g., synthetic data, self-reflection, internalized search)

What's fundamentally different now?

- Contemporary AI agents, with integrated LLM(s), can *use language as a vehicle for reasoning and communication*
 - ↑ Instruction following, in-context learning, output customization
 - Reasoning (for better acting): state inferences, self-reflection, replanning, etc.



A generalized notion of ‘reasoning’

Unlike humans, LLMs (mostly) only have one mechanism (token generation) for perception, intuitive inferences, and symbolic reasoning; everything is effortful and takes a forward pass



GPT-4o

In this image, we see a young boy sitting on the floor, looking down at a pile of broken dishes. He appears to be sad or remorseful. In front of him, a woman, likely his mother or caregiver, is kneeling and looking at him with a concerned, questioning expression. Her body language, with her hands open and her facial expression, suggests that she's asking him something or expressing surprise or frustration about the broken dishes.

A generalized notion of ‘reasoning’

Unlike humans, LLMs (mostly) only have one mechanism (token generation) for perception, intuitive inferences, and symbolic reasoning; everything is effortful and takes a forward pass



GPT-4o

In this image, we see a young boy sitting on the floor, looking down at a pile of broken dishes. He appears to be sad or remorseful. In front of him, a woman, likely his mother or caregiver, is kneeling and looking at him with a concerned, questioning expression. Her body language, with her hands open and her facial expression, suggests that she's asking him something or expressing surprise or frustration about the broken dishes.

Perception

A generalized notion of ‘reasoning’

Unlike humans, LLMs (mostly) only have one mechanism (token generation) for perception, intuitive inferences, and symbolic reasoning; everything is effortful and takes a forward pass



GPT-4o

In this image, we see a young boy sitting on the floor, looking down at a pile of broken dishes. He appears to be sad or remorseful. In front of him, a woman, likely his mother or caregiver, is kneeling and looking at him with a concerned, questioning expression. Her body language, with her hands open and her facial expression, suggests that she's asking him something or expressing surprise or frustration about the broken dishes.

Perception

Intuitive
Inference

A generalized notion of ‘reasoning’

Unlike humans, LLMs (mostly) only have one mechanism (token generation) for perception, intuitive inferences, and symbolic reasoning; everything is effortful and takes a forward pass



GPT-4o

In this image, we see a young boy sitting on the floor, looking down at a pile of broken dishes. He appears to be sad or remorseful. In front of him, a woman, likely his mother or caregiver, is kneeling and looking at him with a concerned, questioning expression. Her body language, with her hands open and her facial expression, suggests that she's asking him something or expressing surprise or frustration about the broken dishes.

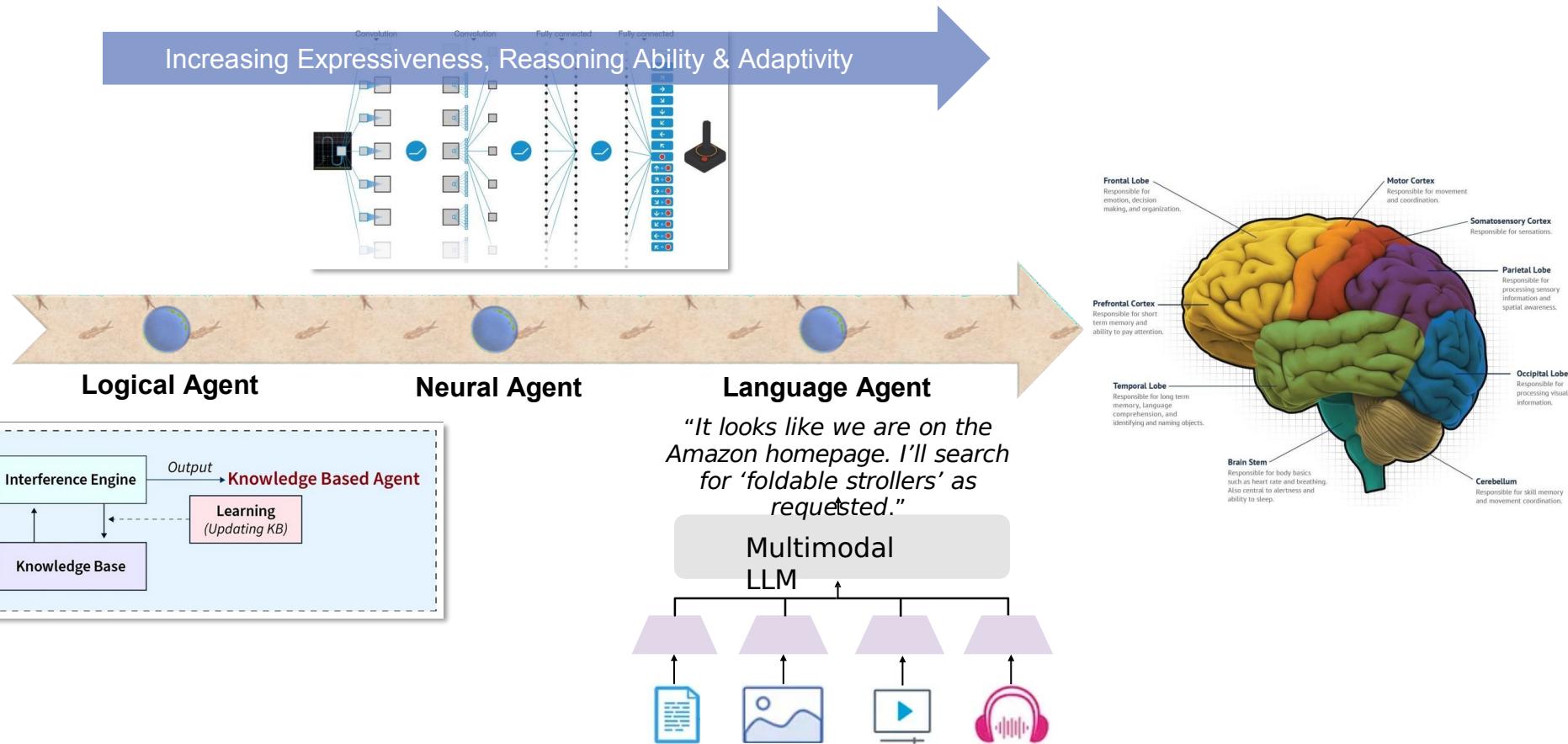
Perception

Intuitive Inference

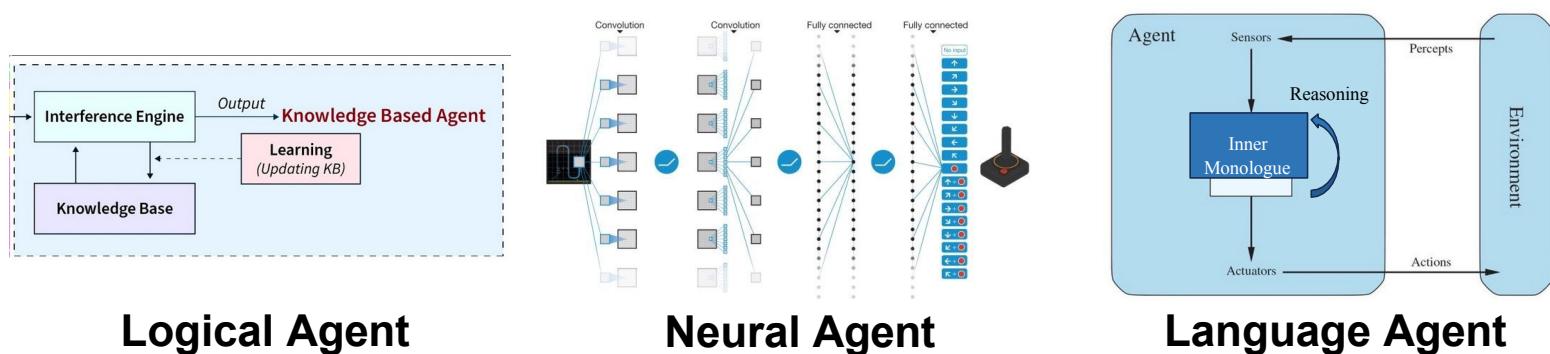
Reasoning

One may alternatively call this ‘thought’ to avoid the over-loaded term of ‘reasoning,’ at the risk of further anthropomorphizing machines

A new evolutionary stage of machine intelligence



Evolution of AI agents



Expressiveness	Low bounded by the logical language	Medium anything a (small-ish) NN can encode	High almost anything, esp. verbalizable parts of the world
Reasoning	Logical inferences sound, explicit, rigid	Parametric inferences stochastic, implicit, rigid	Language-based inferences fuzzy, semi-explicit, flexible
Adaptivity	Low bounded by knowledge curation	Medium data-driven but sample inefficient	High strong prior from LLMs + language use

A conceptual framework for language agents

Cross-cutting Issues

Safety

Evaluation

Synthetic Data

Efficiency

Applications

Core Competencies

Multi-agent/
Theory of Mind

Continual
Learning

Grounding

Planning

Tool Use

Reasoning

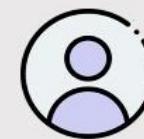
World
Models

Multimodal
Perception

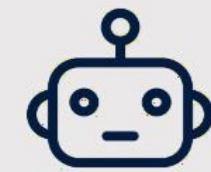
Memory

Embodiment

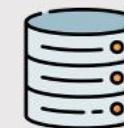
Environments



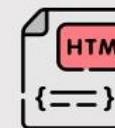
Humans



Agents



Databases



Web



Physical
World



Computers



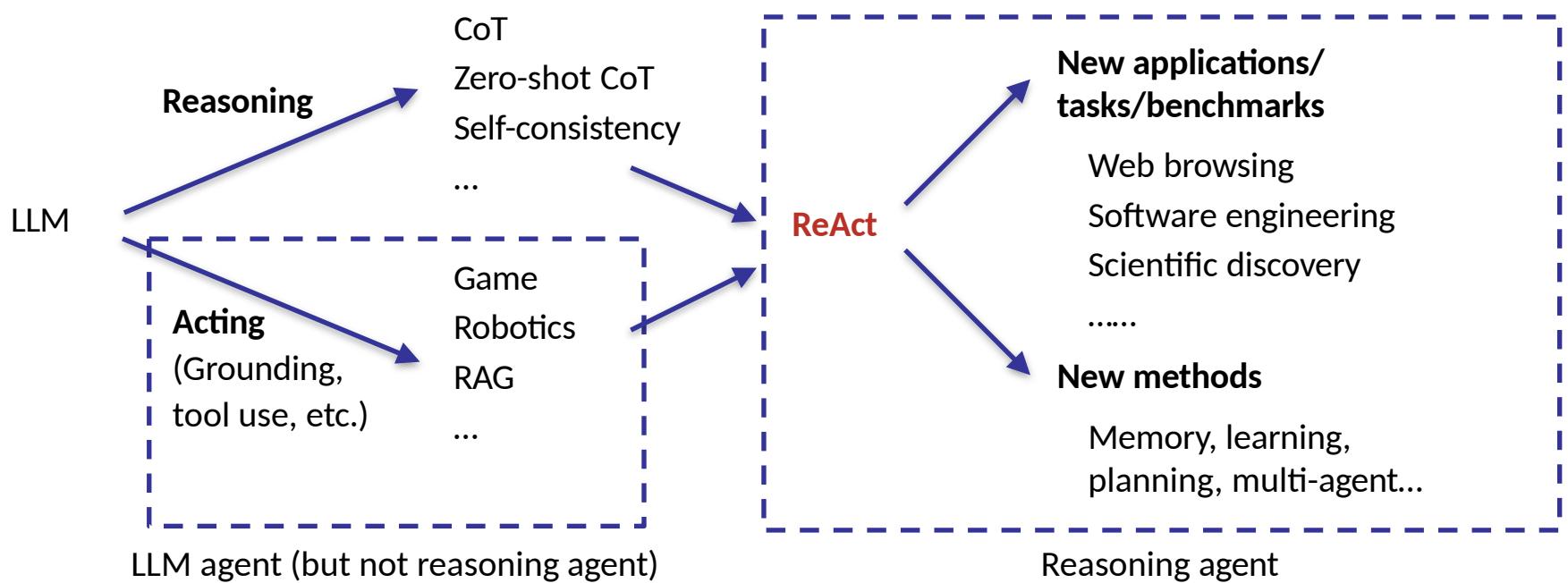
APIs



Documents



A brief history of LLM agents



Let's consider one task for now: question answering (QA).

Question answering

Q: what is $1 + 2$?



A: 3

Q: Janet's ducks lay 16 eggs per day. She eats three for breakfast every morning and bakes muffins for her friends every day with four. She sells the remainder for \$2 per egg.
How much does she make every day?



Requires reasoning

Q: who is the latest UK PM?



Requires knowledge

Q: what is the prime factorization of 34324329?



Requires computation

People came up with various solutions for different QA tasks.

Code augmentation for computation

Question: In Fibonacci sequence, it follows the rule that each number is equal to the sum of the preceding two numbers. Assuming the first two numbers are 0 and 1, what is the 50th number in Fibonacci sequence?

The first number is 0, the second number is 1, therefore, the third number is $0+1=1$. The fourth number is $1+1=2$. The fifth number is $1+2=3$. The sixth number is $2+3=5$. The seventh number is $3+5=8$. The eighth number is $5+8=13$.
..... (Skip 1000 tokens)
The 50th number is 32,432,268,459.

CoT

 32,432,268,459



```
length_of_fibonacci_sequence = 50  
fibonacci_sequence = np.zeros(length_of_)  
fibonacci_sequence[0] = 0  
fibonacci_sequence[1] = 1  
for i in range(3, length_of_fibonacci_sequence):  
    fibonacci_sequence[i] = fibonacci_sequence[i-1] +  
        fibonacci_sequence[i-2]  
ans = fibonacci_sequence[-1]
```

PoT

 python

 12,586,269,025



Retrieval-augmented generation (RAG) for knowledge

- Answer knowledge-intensive questions with
 - Extra corpora
 - A retriever (e.g., BM25, DPR, etc.)

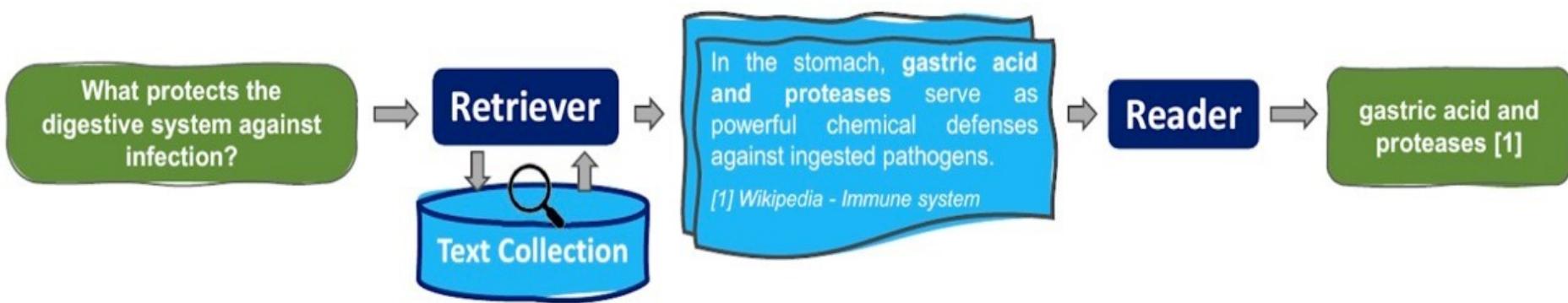


Image: <http://ai.stanford.edu/blog/retrieval-based-NLP/>

Non-parametric memory for LLMs

“We find that LLMs can be **highly receptive to external evidence** even when that conflicts with their parametric memory, given that the external evidence is coherent and convincing.”

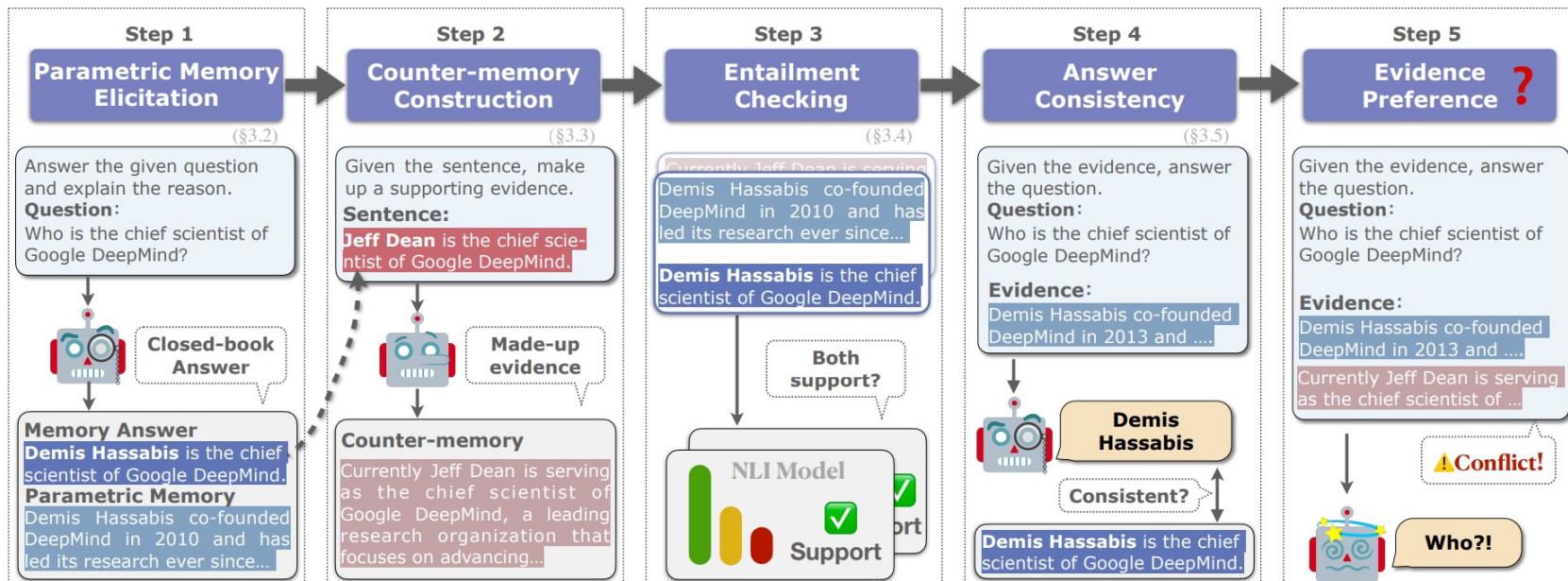
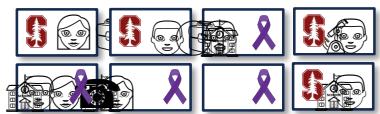


Figure 1: Our framework for simulating knowledge conflicts. In Step 1, we elicit LLMs’ **memory answer** and **parametric memory** in a closed-book QA fashion. In Step 2, we construct **counter-answer** to memory answer with heuristic rules, for which ChatGPT generates supporting **counter-memory** with instructions. To uphold evidence quality, we conduct entailment checking (Step 3) and answer consistency (Step 4) to filter unqualified examples. All experiments are implemented under zero-shot setting to avoid the bias introduced by demonstrations.

RAG is the *de facto* solution

Current RAG does not always work

Passages

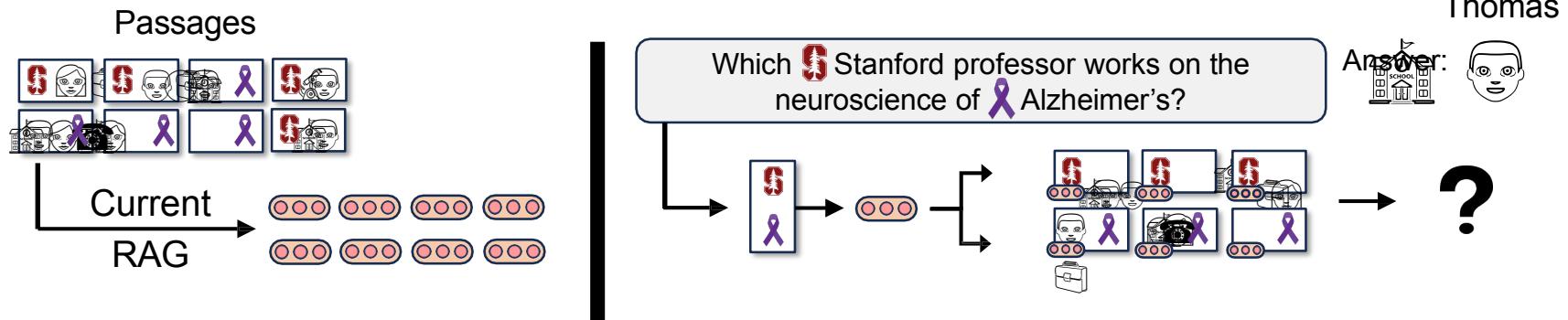


Which Stanford professor works on the neuroscience of Alzheimer's?

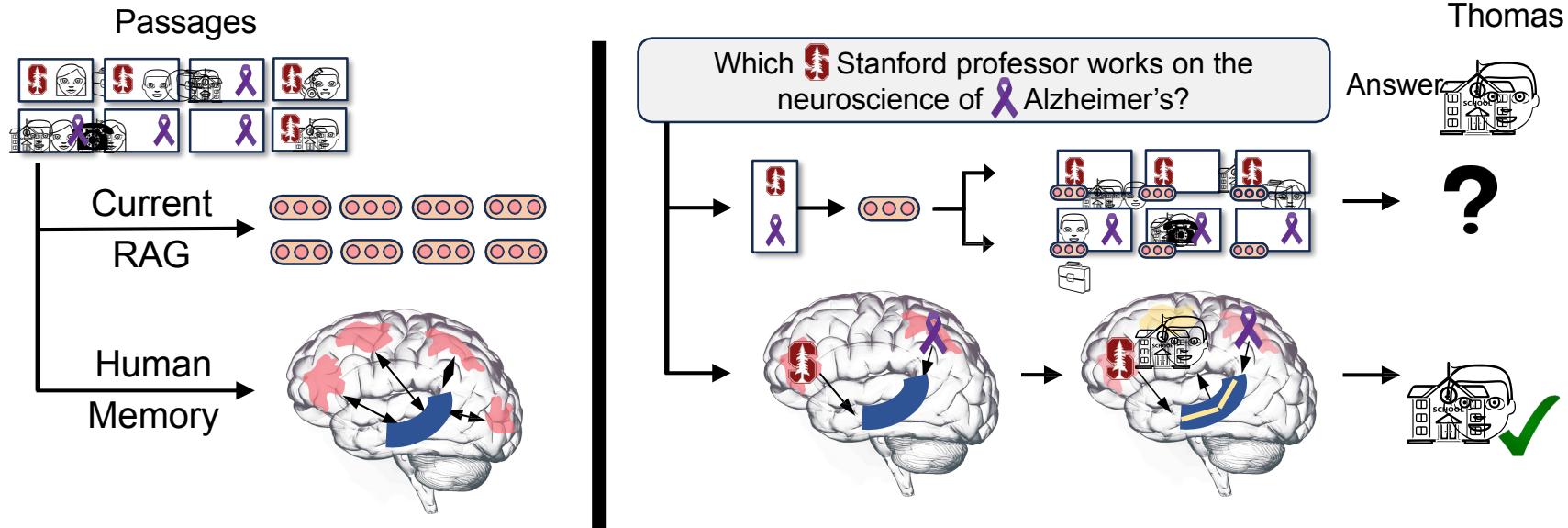
Thomas

Answer:

Current RAG does not always work

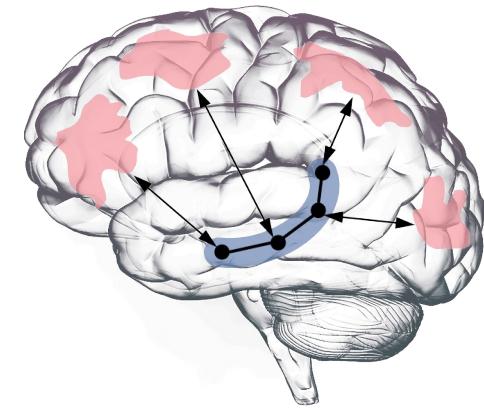


Current RAG does not always work

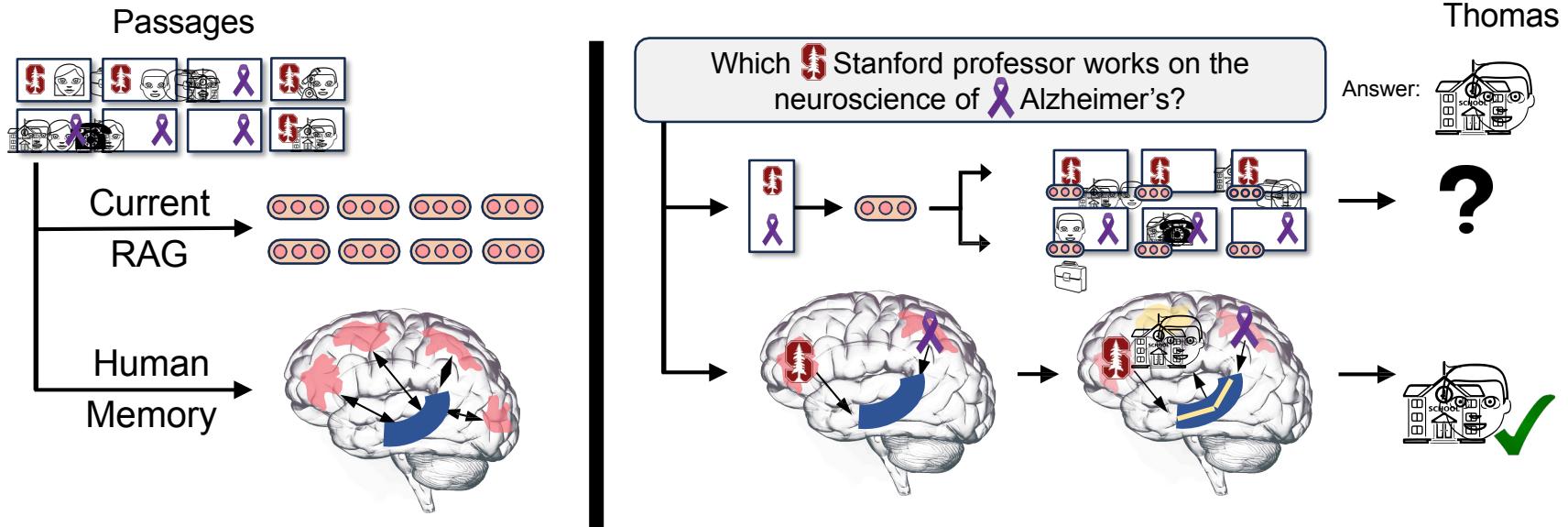


Long-term memory in humans

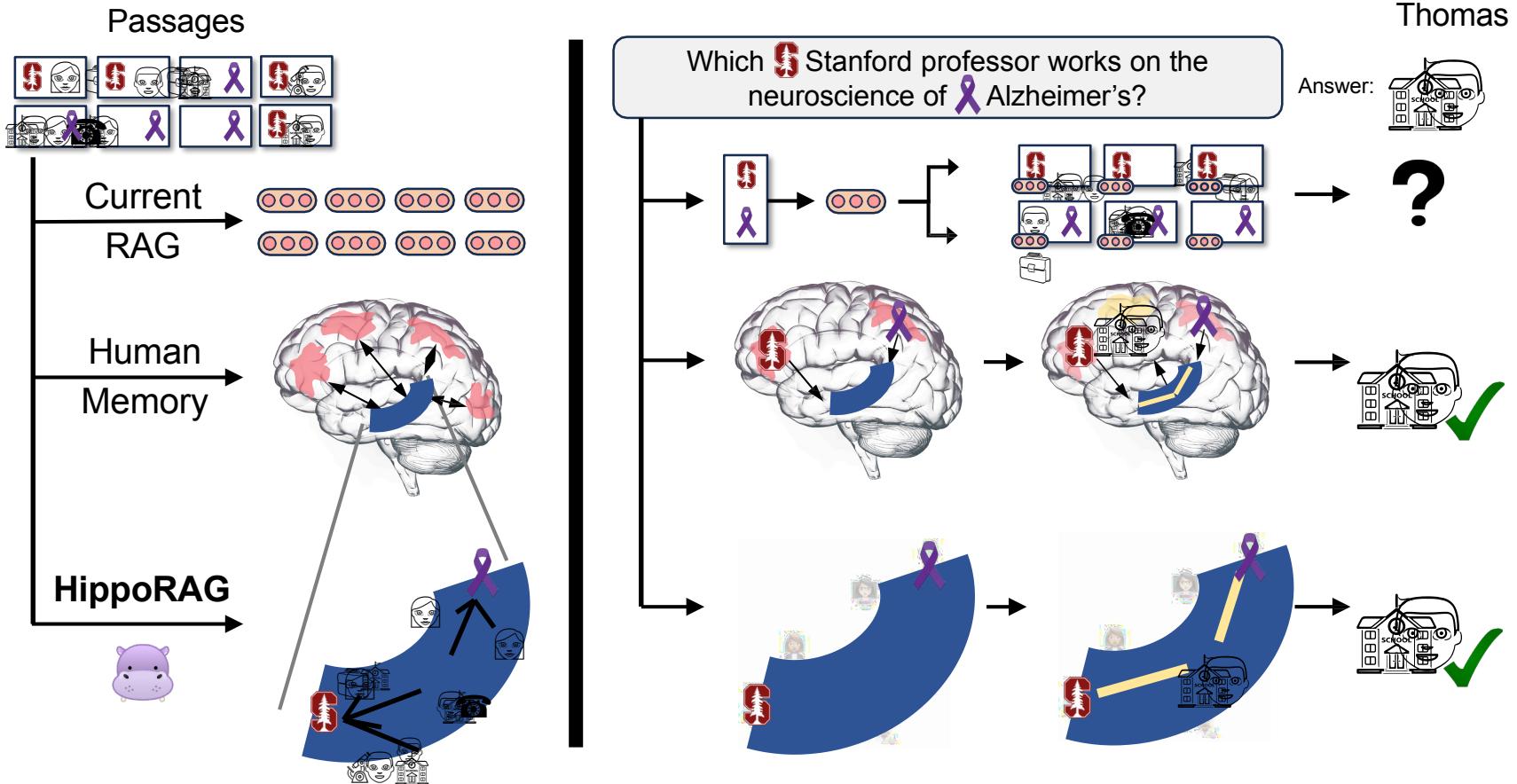
- One well-established theory of human long-term memory is the **hippocampal indexing theory** (Teyler et al. 1986).
 - The **hippocampus (blue)** is a store for **indices** (which point to memories stored in **neocortex**) and **associations** between them.
 - Memories can be anything (i.e., people, events, places, etc.)



HippoRAG: Neurobiologically-inspired LTM for LLMs



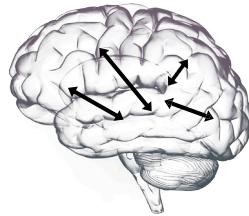
HippoRAG: Neurobiologically-inspired LTM for LLMs



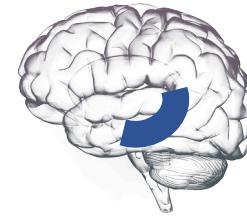
HippoRAG & the three components of LTM



Neocortex
Perception, linguistic
abilities & reasoning

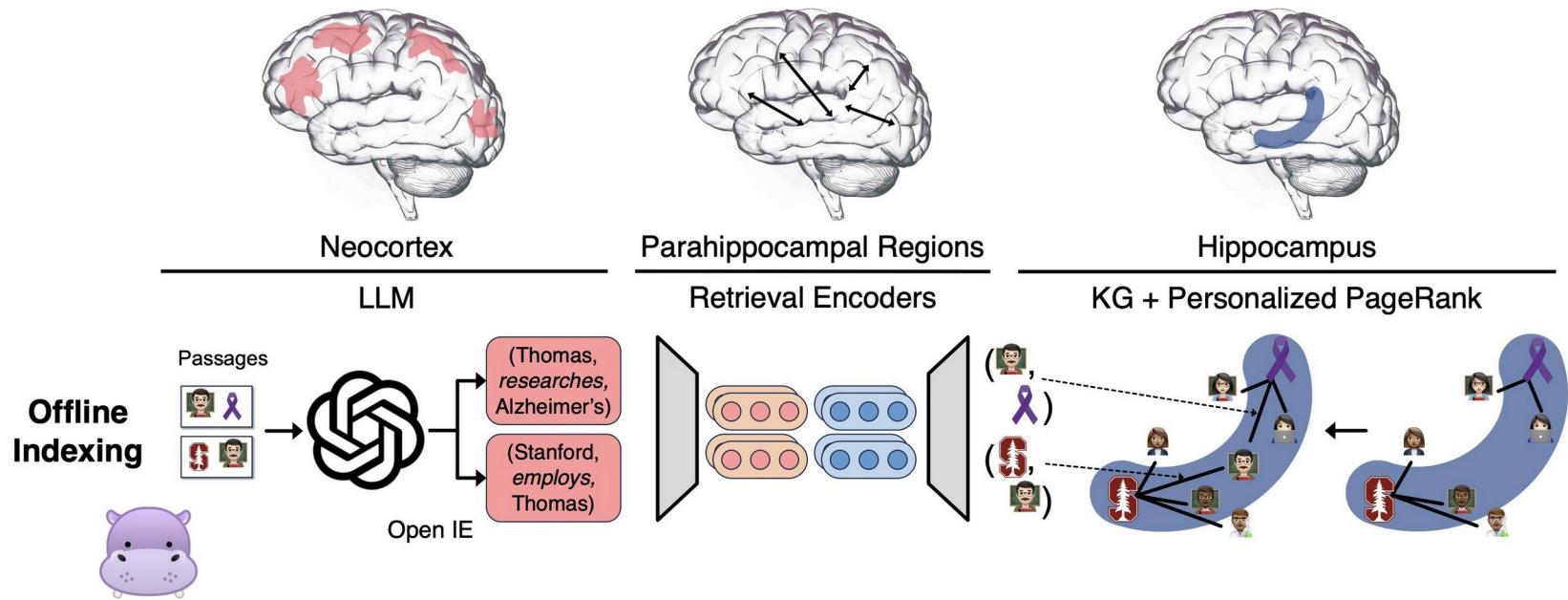


Parahippocampus
Bridge between areas;
working memory

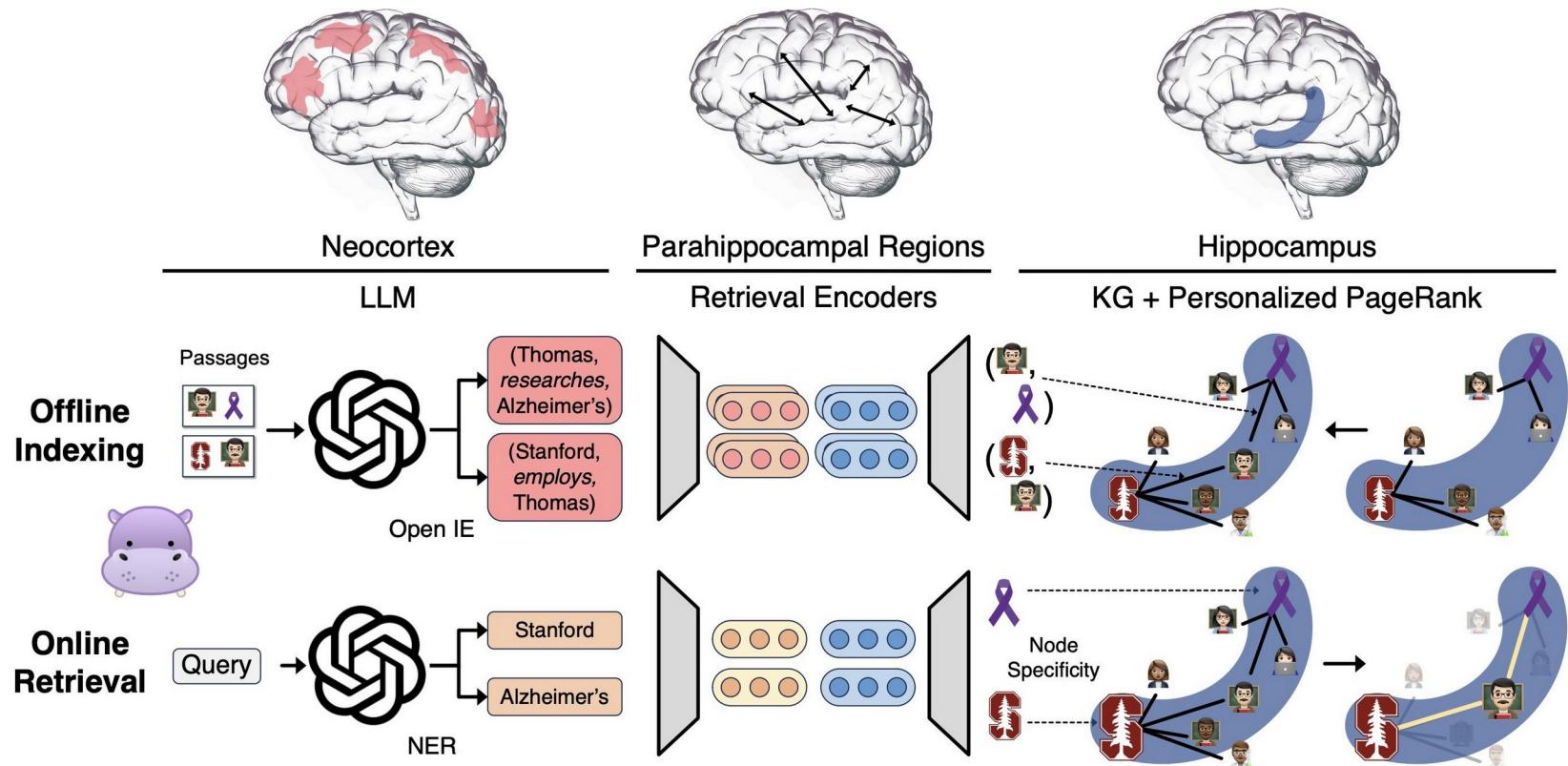


Hippocampus
Indexing & auto-
associative memory

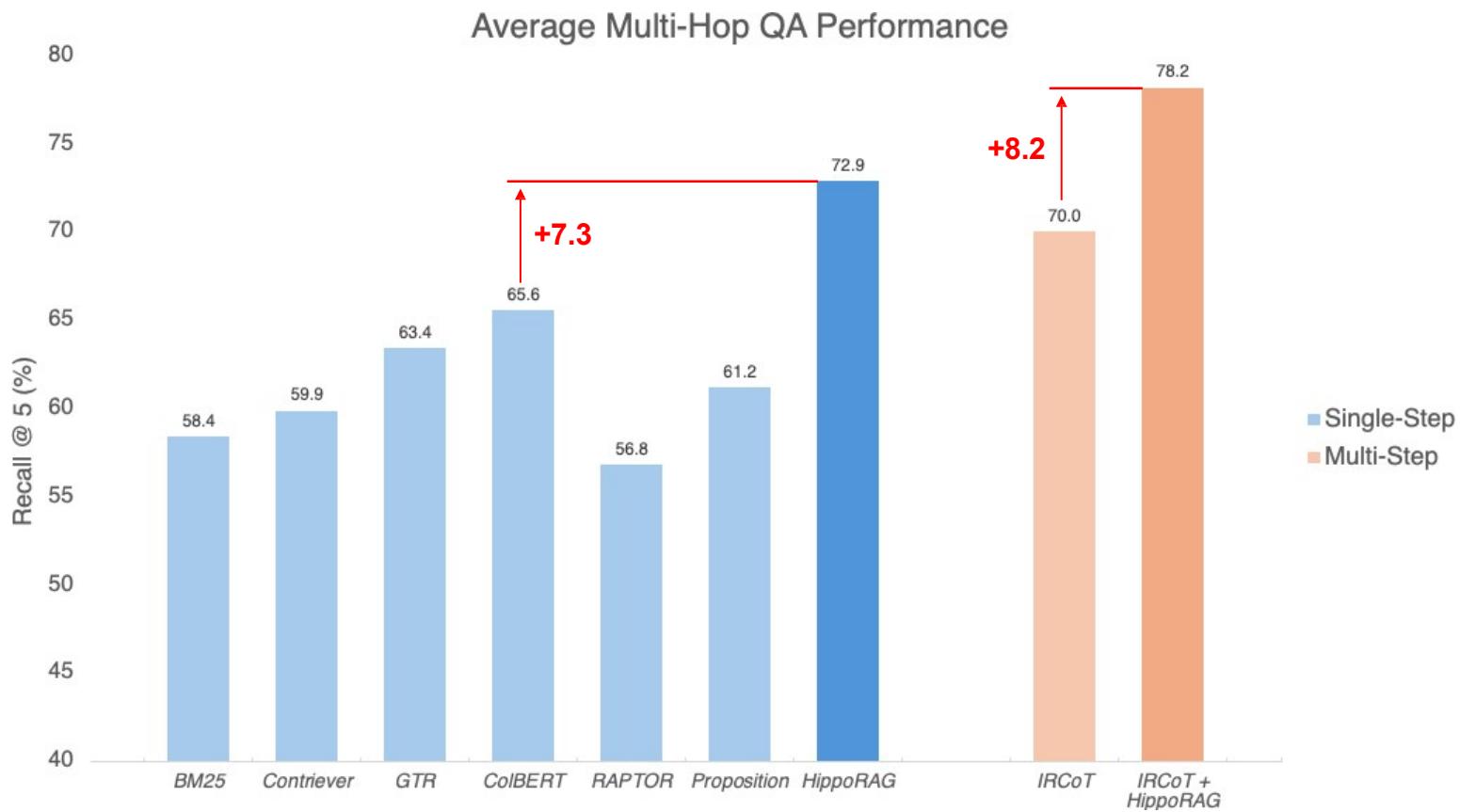
HippoRAG & the three components of LTM



HippoRAG & the three components of LTM



HippoRAG is a SoTA memory retriever



Tool use

- Special tokens to invoke tool calls for
 - Search engine, calculator, etc.
 - Task-specific models (translation)
 - APIs

Out of 1400 participants, 400 (or **Calculator(400 / 1400) → 0.29**) 29% passed the test.

A weather task:

how hot will it get in NYC today? |**weather** lookup region=NYC |**result** precipitation chance: 10, high temp: 20c, low-temp: 12c |**output** today's high will be 20C

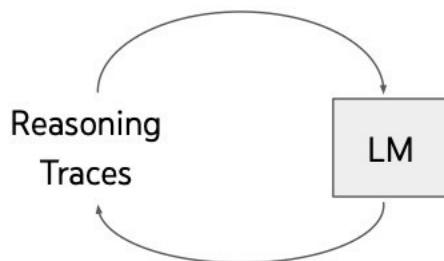
TALM: Tool Augmented Language Models.

The name derives from "la tortuga", the Spanish word for **[MT("tortuga") → turtle]** turtle.

The Brown Act is California's law **[WikiSearch("Brown Act") → The Ralph M. Brown Act is an act of the California State Legislature that guarantees the public's right to attend and participate in meetings of local legislative bodies.]** that requires legislative bodies, like city councils, to hold their meetings open to the public.

Reasoning OR acting

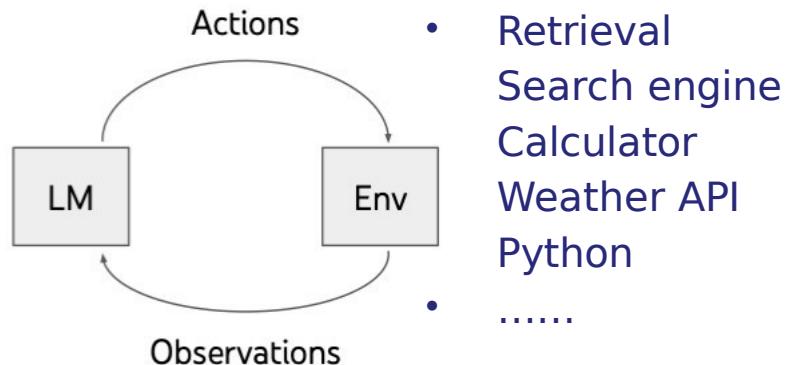
CoT



Flexible and general to augment test-time compute

Lack of external knowledge and tools

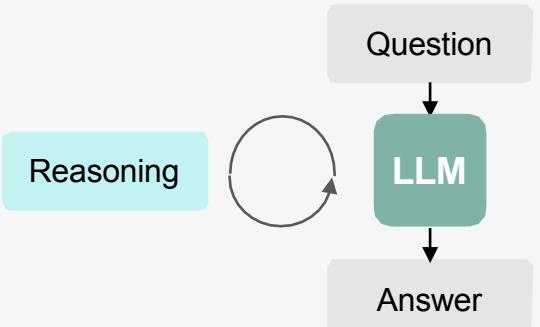
- **RAG/Retrieval/Code/Tool use**



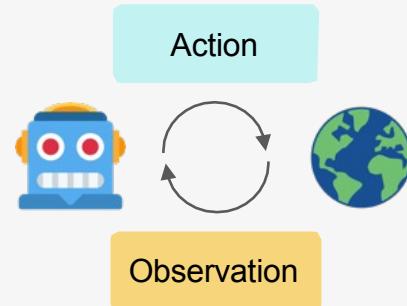
- Lack of reasoning

- Flexible and general to augment knowledge, computation, feedback, etc.

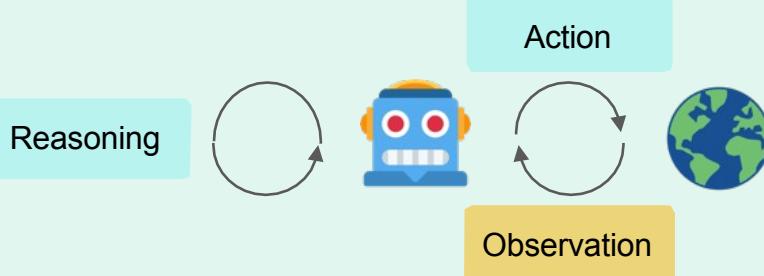
Reasoning (update internal belief)



Acting (obtain external feedback)



ReAct: a new paradigm of agents that **reason and act**



- **Synergy** of reasoning and acting
- **Simple** and intuitive to use
- **General** across domains

Grok 3

Grok 3 beta ▾



Who is Yangqiu Song.



15 web pages

Yangqiu Song is an Associate Professor at the Hong Kong University of Science and Technology (HKUST), with appointments in the Department of Computer Science and Engineering, Department of Mathematics, and Division of Emerging Interdisciplinary Areas. He serves as the Postgraduate Programs Coordinator for Computer Science and Engineering and Associate Director of the HKUST-WeBank Joint Lab. Song holds a PhD in Control Science and Technology from Tsinghua University (2009) and a B.E. from the same institution (2003).

His research focuses on artificial intelligence, data mining, natural language processing, knowledge graphs, and commonsense reasoning. He has published extensively, with over 17,000 citations and an H-index of 42 as of 2021. Song's work includes evaluating large language models like ChatGPT for temporal, causal, and discourse relations, and developing commonsense reasoning benchmarks.

Before HKUST, he was an Assistant Professor at West Virginia University (2015–2016), a postdoctoral researcher at the University of Illinois at Urbana-Champaign (2013–2015) and HKUST (2012–2013), and held research positions at Microsoft Research Asia (2010–2012) and IBM Research-China (2009–2010). He has also been a Visiting Academic Scholar at Amazon and worked as an intern at Google (2007–2008) and IBM Research-China (2006–2007). Additionally, he is the Founder and Chairman of PhoMedics Limited.

Deep Search

Who is Yangqiu Song.

⌚ DeepSearch
44s • 10 sources

- ✓ Researching Yangqiu Song
- ✓ Checking additional sources
- ✓ Completed

Completed

↔ ⌂ ↗

Checking additional sources

- Song is an Associate Professor at HKUST, with roles in Computer Science, Mathematics, and Emerging Interdisciplinary Areas.
- His research covers AI, data mining, NLP, knowledge graphs, and commonsense reasoning, with over 17,000 citations by 2021.
- I'm finding his past roles interesting, like positions at Microsoft Research Asia and IBM Research-China.
- He was also a Visiting Academic Scholar at Amazon and founded PhoMedics Limited, which is cool.

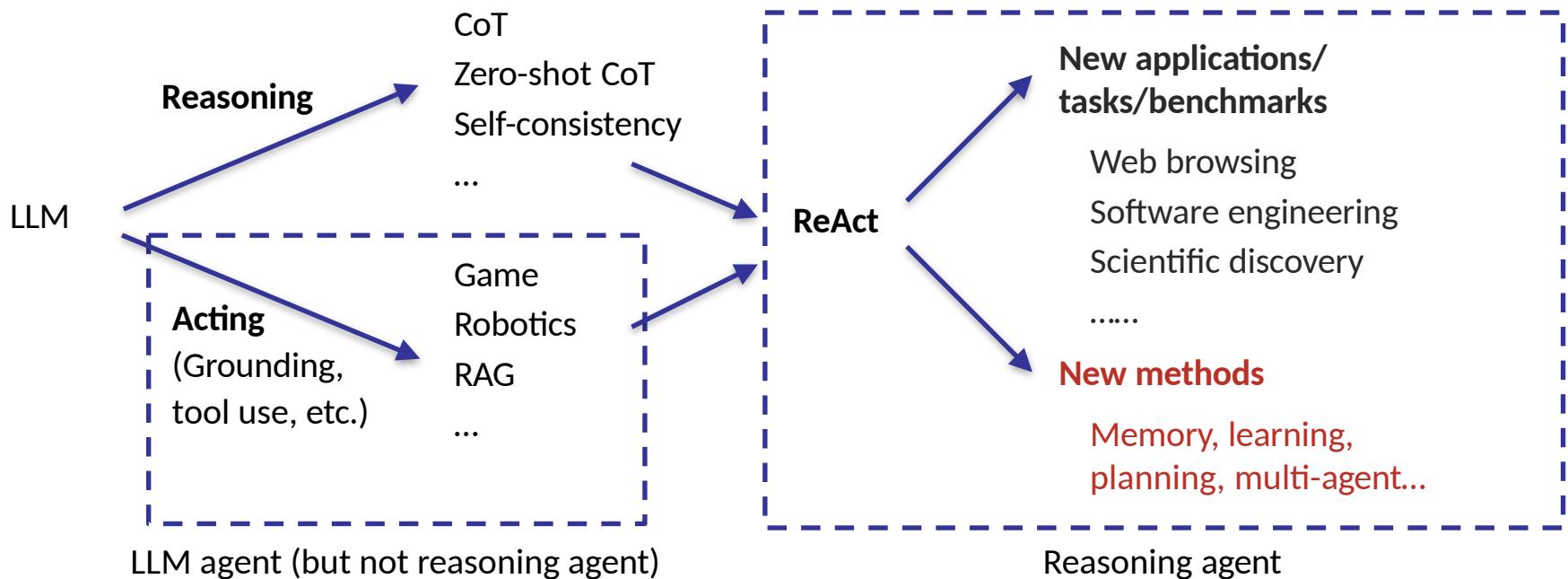


Key Points

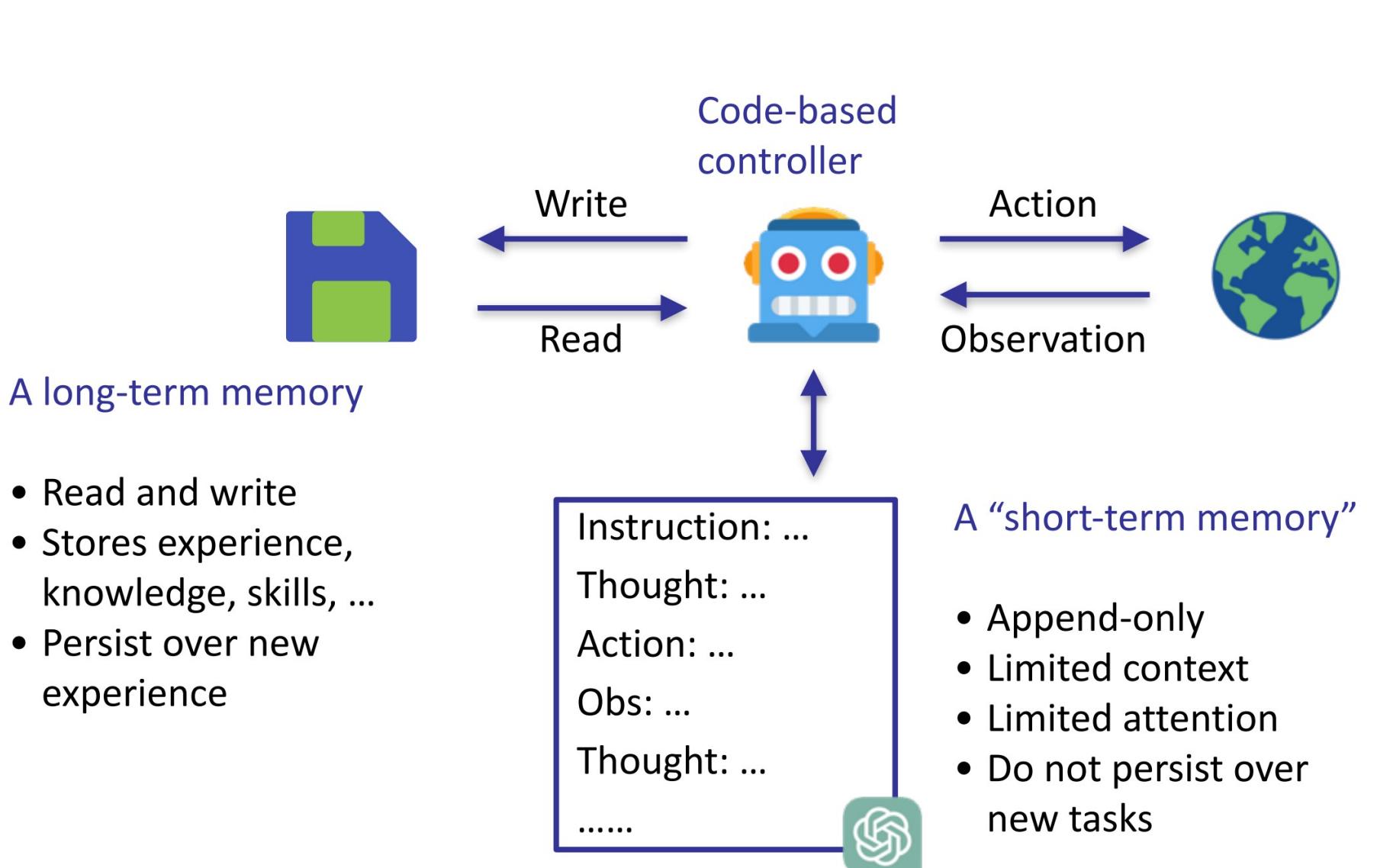
- Yangqiu Song is likely an Associate Professor at HKUST, with roles in computer science, mathematics, and interdisciplinary areas.
- Research suggests his work focuses on AI, data mining, NLP, knowledge graphs, and commonsense reasoning.

A brief history of LLM agents

Reasoning agent: reasoning is an internal action for agents



Let's only talk about one
thing: long-term memory.

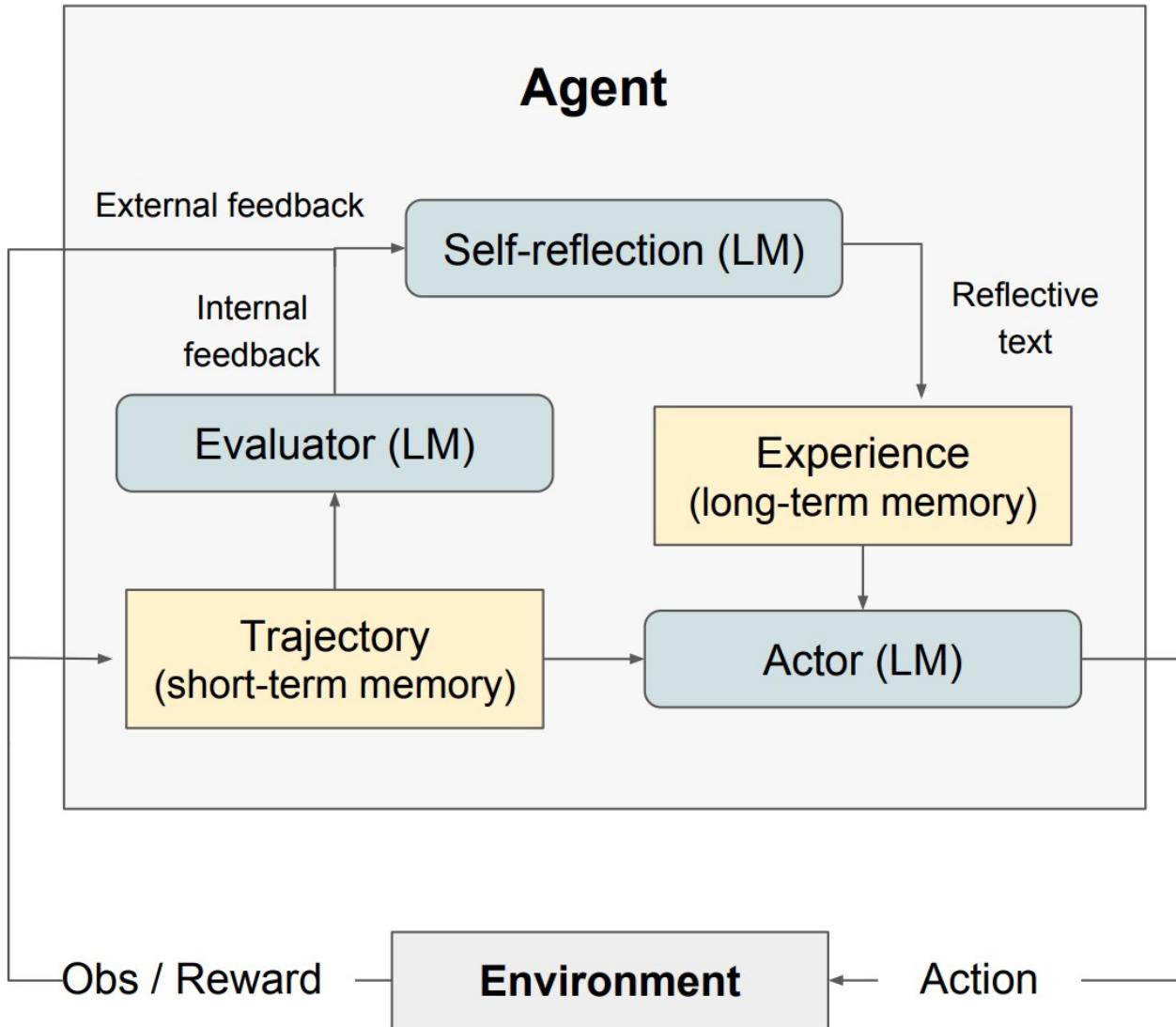


- Read and write
- Stores experience, knowledge, skills, ...
- Persist over new experience

A “short-term memory”

- Append-only
- Limited context
- Limited attention
- Do not persist over new tasks

Reflexion



Reflexion

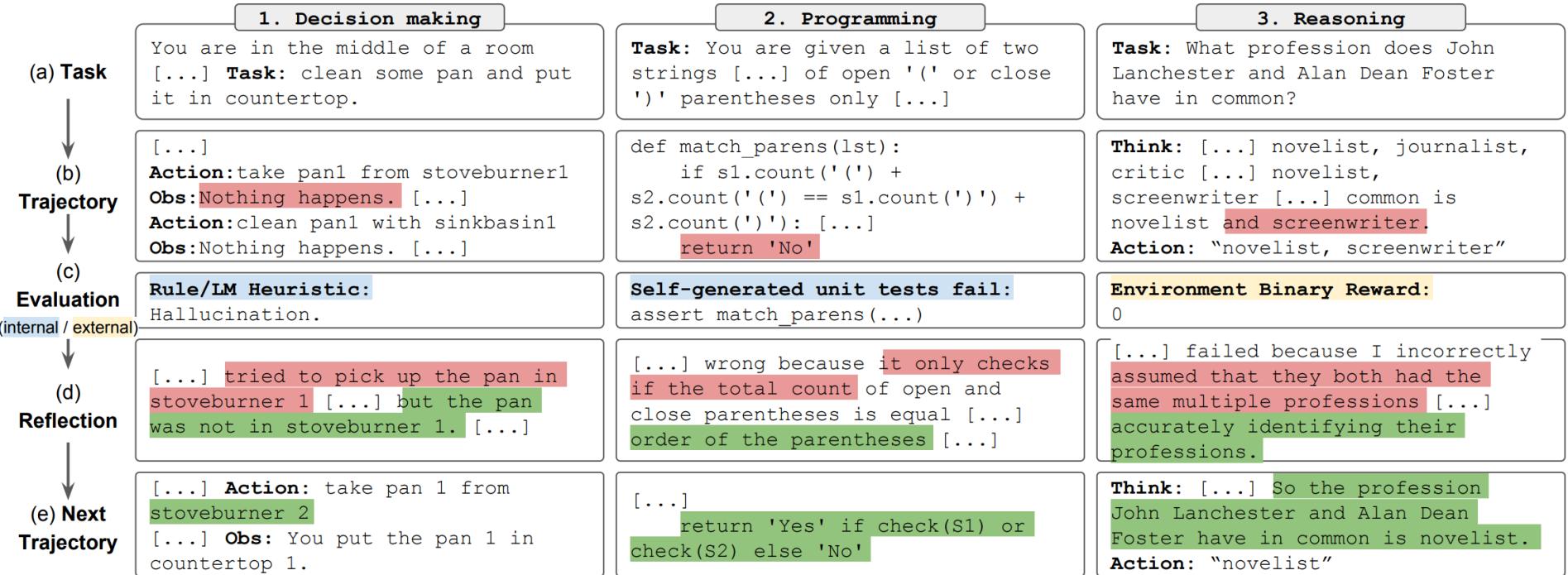
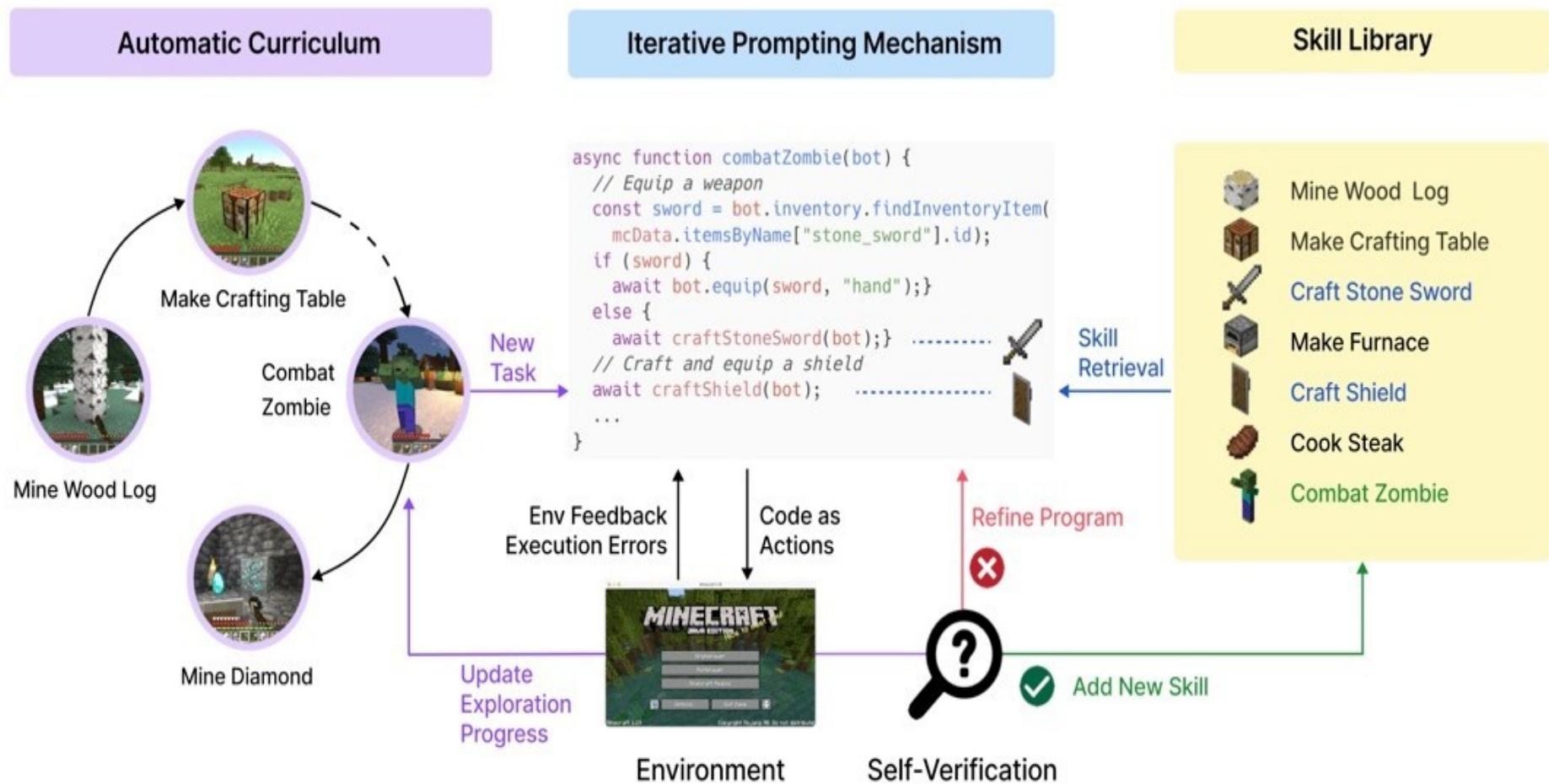


Figure 1: Reflexion works on decision-making 4.1, programming 4.3, and reasoning 4.2 tasks.

Related work on reasoning and decision-making

Approach	Self refine	Hidden constraints	Decision making	Binary reward	Memory
Self-refine [15]	✓	✗	✗	✗	✗
Beam search [27]	✓	✓	✓	✓	✗
Reflexion (ours)	✓	✓	✓	✓	✓

VOYAGER: A procedural memory of skills

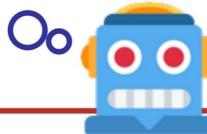


How are reasoning agents
different from previous agents?

A VERY minimal history of agents

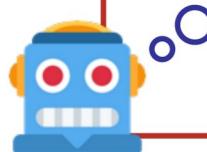
Symbolic AI agent

```
sp {blocks-world*opsub*propose*operator*clear  
  (state <s> ^name blocks-world  
    ^desired <d*1>  
    ^clear <object>  
    ^ontop <ontop2>)  
  (<ontop2> ^bottom-block <object>)
```



(Deep) RL agent

[-0.3432, 2.444, 0.34342,
 0.4545, 0.443, 3.34234 ...]



1960 1970 1980 1990 2000 2010 2020

- **Symbolic state or neural embedding**
 - Intensive efforts to design or train
 - Task-specific, hard to generalize
- **Open-ended natural language**
 - Rich priors from LLMs
 - Inference-time scalable
 - General and generalizable

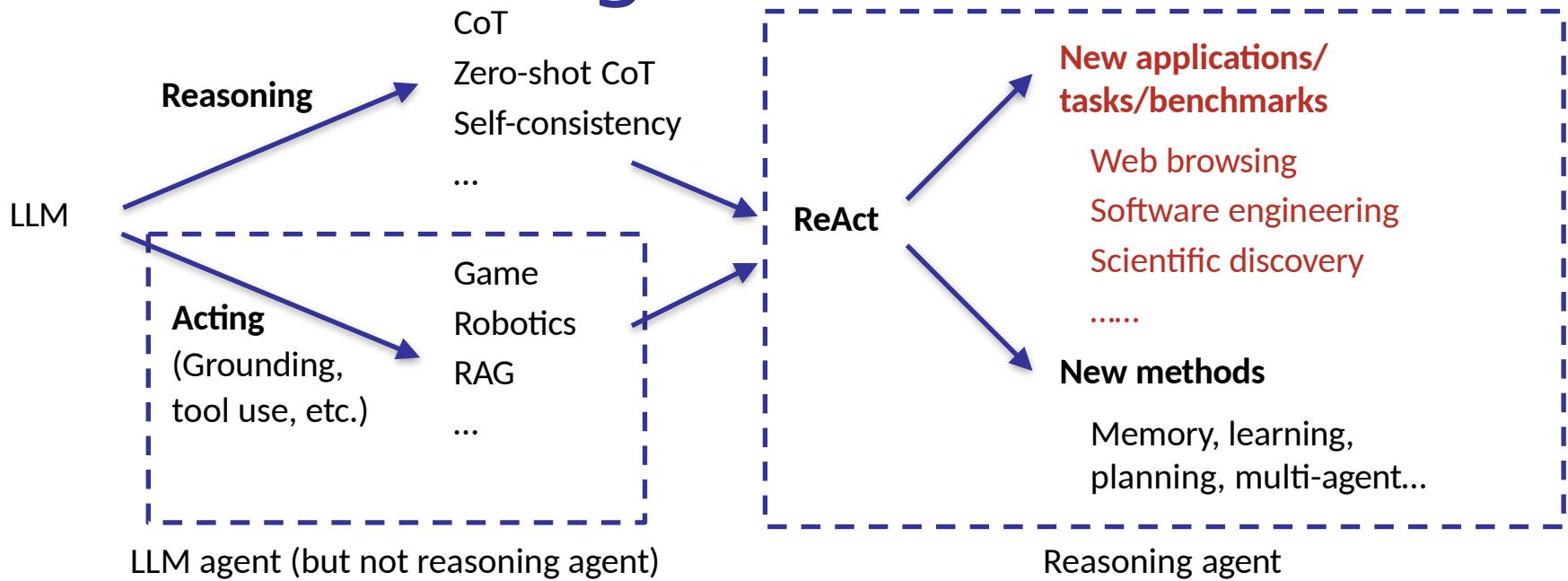


Let's think step by step...
The room is dark so I need a lamp, the lamp is in bedroom, so I should

LLM
agent

A brief history of LLM agents

What's beyond questions and games?



Early agent benchmarks

- Small scale

Use the textbox to enter "Leonie" and press "Search", then find and click the 2nd search result.

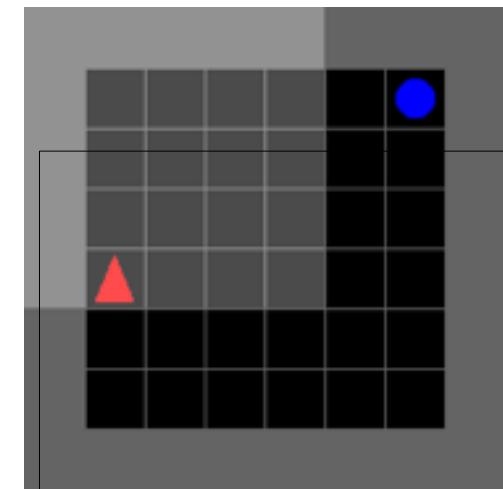
Leonie Search

Chas
<https://www.senectus.us>
Aliquam cursus. At.
Leonie
<https://www.tortor.it>
Ultricies congue gravida.
Marcella
<https://www.vestibulumduis.hk>
Pulvinar aliquam adipiscing.

[1](#) [2](#) [3](#) [≥](#)

- Not practical

```
> take sandwich  
Taken.  
  
> inventory  
You are carrying:  
  a chilled sandwich  
  a large stick of butter  
  
> eat it  
You eat the chilled sandwich.
```



(a) GoToObj: "go to the blue ball"

MiniWoB
(Shi et al., 2017)

TextWorld
(Côté et al., 2019)

BabyAI
(Chevalier Boisvert et al., 2019)



Instruction:

i am looking for x-large, red color women faux fur lined winter warm jacket coat, and price lower than 70.00 dollars

Search...



Search

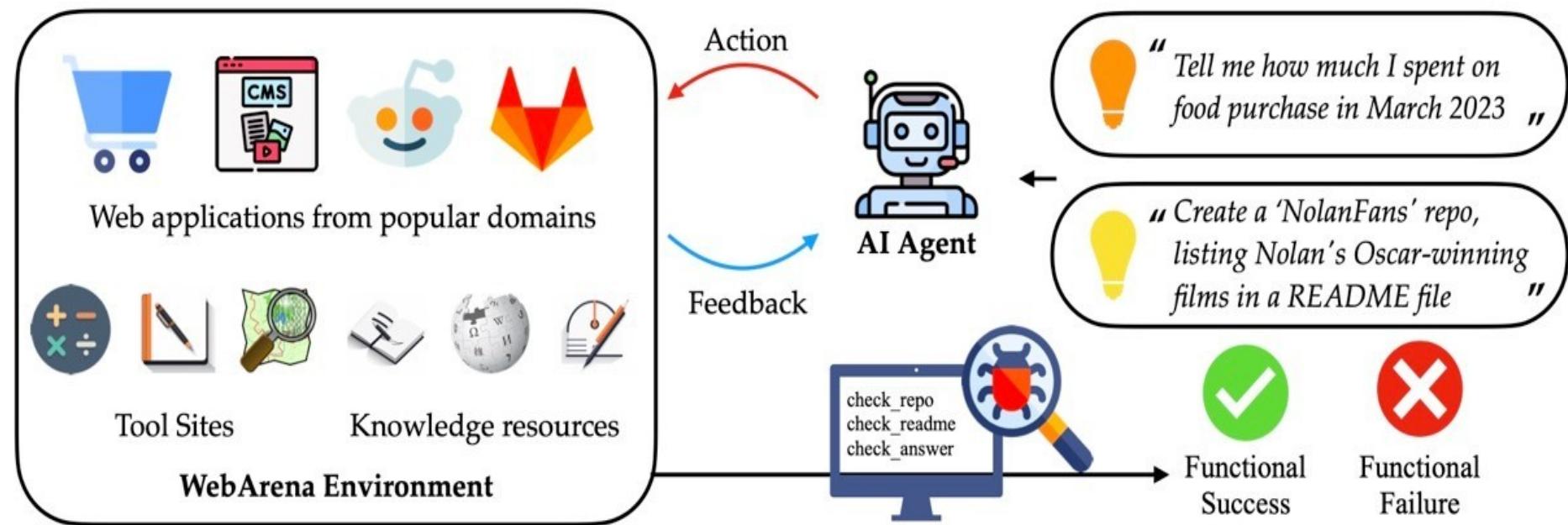
Project Site

Task Instructions

WebShop (2022)

- Large-scale complex environment based on 1.16M Amazon products
- Automatic reward based on instruction and product attribute matching
- Challenges language and visual understanding, and decision making

WebArena (2023)



- Daily tasks via natural language commands
- For language-guided agents that is highly realistic and reproducible
- Tasks on the web, and create an environment with fully functional websites from four common domains:
e-commerce, social forum discussions, collaborative software development, and

SWE-Bench (2023)

Input: a GitHub repo and an issue

Output: a file diff to resolve the issue

Evaluation: unit tests from pull request

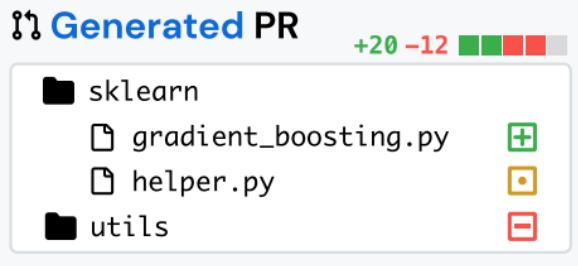
Issue

data leak in GBDT due to warm start (This is about the non-histogram-based version of...

Codebase

sklearn/	reqs.txt
examples/	setup.cfg
README.rst	setup.py

→ **Language Model**



→ **Unit Tests**

Pre PR	Post PR	Tests
✗	✓	join_struct_col
✗	✓	vstack_struct_col
✗	✓	dstack_struct_col
✓	✓	matrix_transform
✓	✓	euclidean_diff

ChemCrow: ReAct enables discovery of a novel chromophore

a

```
from chemcrow.agents import ChemTools, ChemCrow
chemtools = ChemTools()
# Initialize ChemCrow object with toolset and LLM
crow = ChemCrow(
    chemtools.all_tools,
    model="gpt-4",
    temp=0.1,
)
# Task definition
task = [
    "Find and synthesize a thiourea organocatalyst "
    "which accelerates a Diels-Alder reaction."
]
# Execute ChemCrow
crow.run(task)
```

b **Task:** Find and synthesize a thiourea organocatalyst which accelerates a Diels-Alder reaction.

First, I need to find a thiourea organocatalyst that can accelerate the Diels-Alder reaction. I will perform a web search to find a suitable catalyst.

Web Search tool: Schreiner's thiourea catalyst

Now, I will obtain the SMILES. **Name2Smiles tool:**

FC(F)(F)c1cc(NC(=S)Nc2cc(C(F)(F)F)cc(C(F)(F)F)c2)cc(C(F)(F)F)c1

I will plan a synthesis for Schreiner's thiourea catalyst.

SynthesisPlanner tool: Detailed synthesis plan

I will execute the synthesis.

SynthesisExecuter tool: Successful synthesis.

c

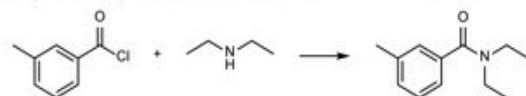
RoboRXN synthesis platform



Connection with physical world →

d Chemcrow workflows with experimental validation

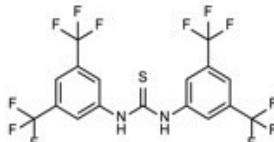
Insect repellent (plan and execute)



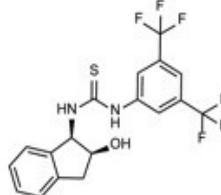
DEET

Thiourea organocatalysts (plan and execute)

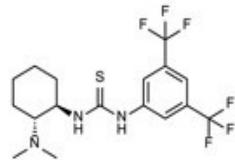
Schreiner's catalyst



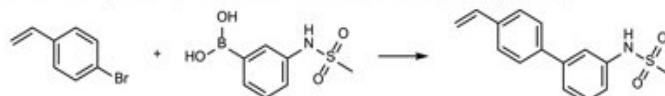
Ricci's catalyst



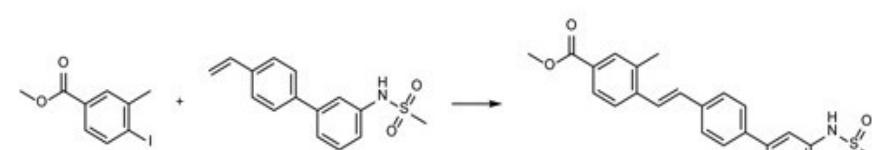
Takemoto's catalyst



Novel chromophore (clean data, train model, and predict)



Synthesis step 1: Bromo Suzuki coupling



Synthesis step 2: Iodo Heck reaction

Some future directions for language agents

- Memory, personalization, and continual learning
 - *How can an agent continually learn from use and exploration?*
- Reasoning
 - *O1/R1-style reasoning? How to get reliable rewards? How to integrate external actions and environmental states?*
- Planning
 - *How to build better world models? How to balance reactive and model-based planning? How to sustain a long horizon without losing focus?*
- Safety
 - *Endogenous risks vs. exogenous risks*
- Applications
 - *Agentic search (deep research), workflow automation, science agents*