



# CS 639: Foundation Models **Agents I**

Fred Sala

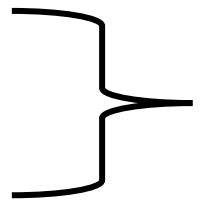
University of Wisconsin-Madison  
**April 21, 2026**



# Announcements

- **Homework 4**: due
- **Project**---ongoing!
- **Next 2 classes (4/23 + 4/28)**: will happen offline; I will release lecture recordings
  - I will be in Brazil for ICLR
- **Class outline**

Tuesday April 21	Agents I
Thursday April 23	Agents II
Tuesday April 28	Applications
Thursday April 30	Future Areas



End of Class

# Outline

- **Introduction to LLM-Powered Agents**

- Motivation, goals, differences vs classical agents / standard LMs, overall architectures and key components

- **Multiagent Systems**

- Motivation, simulations, architectures and design, communication

- **Challenges**

- Scaling, realism, open problems

# Outline

- **Introduction to LLM-Powered Agents**

- Motivation, goals, differences vs classical agents / standard LMs, overall architectures and key components

- **Multiagent Systems**

- Motivation, simulations, architectures and design, communication

- **Challenges**

- Scaling, realism, open problems

# Motivation: From LLMs To Agents

Standard LLMs are static

- But we want **interactive** systems
  - Want systems that can reason, use tools, and act autonomously
  - Want FMs that can be actors, not just predictors
- Useful survey: Wang et al '25  
"Large Language Model Agent: A Survey on Methodology"

## Large Language Model Agent: A Survey on Methodology, Applications and Challenges

Junyu Luo, Weizhi Zhang, Ye Yuan, Yusheng Zhao, Junwei Yang, Yiyang Gu, Bohan Wu, Binqi Chen, Ziyue Qiao, Qingqing Long, Rongcheng Tu, Xiao Luo, Wei Ju, Zhiping Xiao, Yifan Wang, Meng Xiao, Chenwu Liu, Jingyang Yuan, Shichang Zhang, Yiqiao Jin, Fan Zhang, Xian Wu, Hanqing Zhao, Dacheng Tao, *Fellow, IEEE*, Philip S. Yu, *Fellow, IEEE* and Ming Zhang

**Abstract**—The era of intelligent agents is upon us, driven by revolutionary advancements in large language models. Large Language Model (LLM) agents, with goal-driven behaviors and dynamic adaptation capabilities, potentially represent a critical pathway toward artificial general intelligence. This survey systematically deconstructs LLM agent systems through a methodology-centered taxonomy, linking architectural foundations, collaboration mechanisms, and evolutionary pathways. We unify fragmented research threads by revealing fundamental connections between agent design principles and their emergent behaviors in complex environments. Our work provides a unified architectural perspective, examining how agents are constructed, how they collaborate, and how they evolve over time, while also addressing evaluation methodologies, tool applications, practical challenges, and diverse application domains. By surveying the latest developments in this rapidly evolving field, we offer researchers a structured taxonomy for understanding LLM agents and identify promising directions for future research. The collection is available at <https://github.com/luo-junyu/Awesome-Agent-Papers>.

**Index Terms**—Large language model, LLM agent, AI agent, intelligent agent, multi-agent system, LLM, literature survey

### 1 INTRODUCTION

Artificial Intelligence is entering a pivotal era with the emergence of LLM agents—intelligent entities powered by large language models (LLMs) capable of perceiving environments, reasoning about goals, and executing actions [1]. Unlike traditional AI systems that merely respond to user inputs, modern LLM agents actively engage with their environments through continuous learning, reasoning, and adaptation. This shift represents a technological advancement and a fundamental reimagining of human-machine relationships. Commercial LLM agent systems (*e.g.*, DeepResearch, DeepSearch, and Manus) exemplify this

once required human expertise, from in-depth research to computer operation, while adapting to specific user needs.

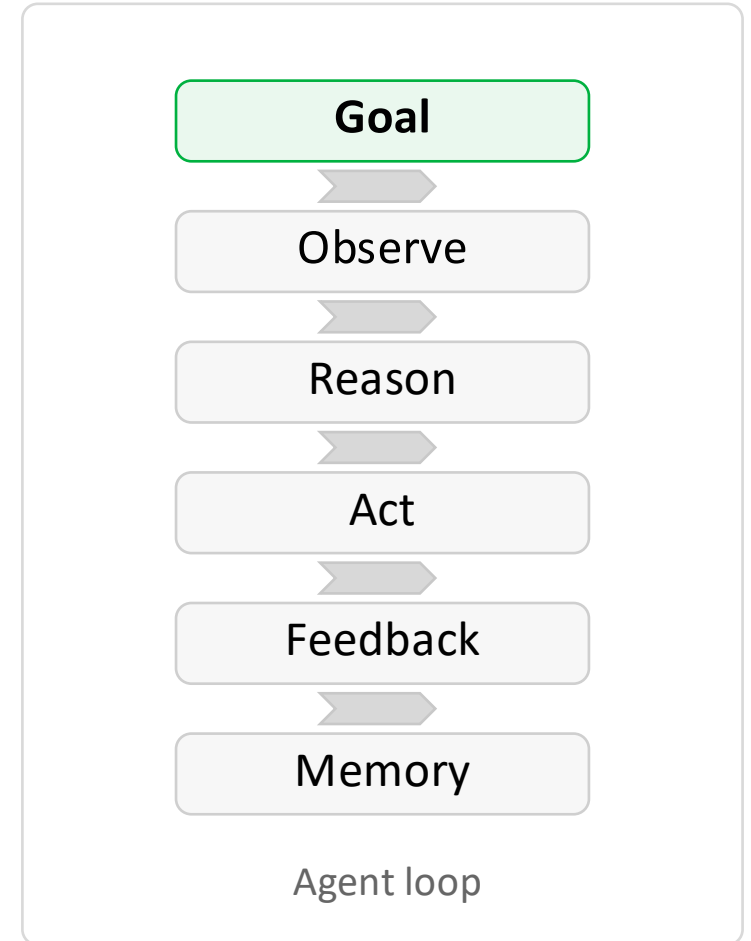
Compared to traditional agent systems [2], LLM-based agents have achieved generational across multiple dimensions, including knowledge sources [3], generalization capabilities [4], and interaction modalities [5]. Today's agents represent a qualitative leap driven by the convergence of three key developments: ① unprecedented reasoning capabilities of LLMs [6], ② advancements in tool manipulation and environmental interaction [7], and ③ sophisticated memory architectures that support longitudinal experience accumulation [8], [9]. This convergence has transformed theoretical constructs into practical systems, increasingly blurring the

503.21460v1 [cs.CL] 27 Mar 2025

# What's an Agent?

Let's **differentiate** from a standard LLM

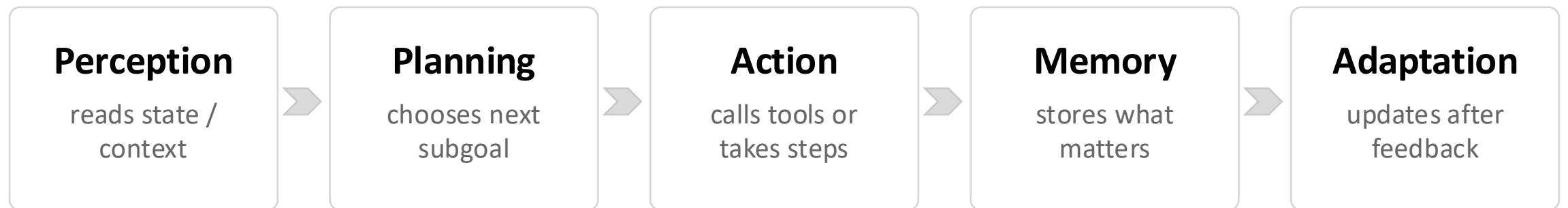
- These usually answer a prompt (single response)
- Agents act, often taking many steps, to accomplish a **goal**
  - Interact with tools, environments, feedback loops
- But, still **built on top of LLMs**
- What do we use them for?
  - As we saw in benchmarks, coding, computer use, research, science, etc.



# More Concrete Definition

An LLM agent is a system that uses a FM to **reason, plan, and act** in an environment

- Often with tools and memory.
- Note: the agent includes the **harness/scaffold**, not just the FM or LLM that powers it



# Pre-LLMs: Classical Agents

Much of the work on agents predates emergence of FMs and LLMs

- A common framework: **PEAS** (Performance, Environment, Actuators, Sensors)
- Agent types:
  - Reactive, deliberative, hybrid
- **Reinforcement learning** is a standard paradigm for agents
  - Convenient since we already use it for LLMs



# Architecture of LLM-Powered Agents

High-level components: **input**, **memory**, **planner**, **tools**,

- **Input**: whatever we have access to
- **Memory**: as simple as context window, but
- **Planning**: reasoning-based decomposition etc
- **Tools**: similar to our earlier discussion on integrating tool use



# Single-Agent vs. Multi-Agent

We can create multiple agents

- These get to interact with each other to
- Note that we can simulate this with a single agent too!
- Key difference: multiple agents permits using **different** harnesses and models
  - This specialization can be more efficient: use cheaper models for easier tasks, etc.
  - But...more complexity!



# Single Agent Patterns

Variations based on where planning, execution, and verification are placed. Some examples:

- **ReAct**
  - Interleave reason and execution
- **Plan/execute**
  - Planning phase followed by execution
- **Reflection**
  - Generate, critique, revise
- **Verifier-based**

REACT: SYNERGIZING REASONING AND ACTING IN LANGUAGE MODELS

Shunyu Yao<sup>\*1</sup>, Jeffrey Zhao<sup>2</sup>, Dian Yu<sup>2</sup>, Nan Du<sup>2</sup>, Izhak Shafran<sup>2</sup>, Karthik Narasimhan<sup>1</sup>, Yuan Cao<sup>2</sup>

<sup>1</sup>Department of Computer Science, Princeton University

<sup>2</sup>Google Research, Brain team

<sup>1</sup>{shunyuy, karthikn}@princeton.edu

<sup>2</sup>{jeffreyzhao, dianyu, dunan, izhak, yuancao}@google.com

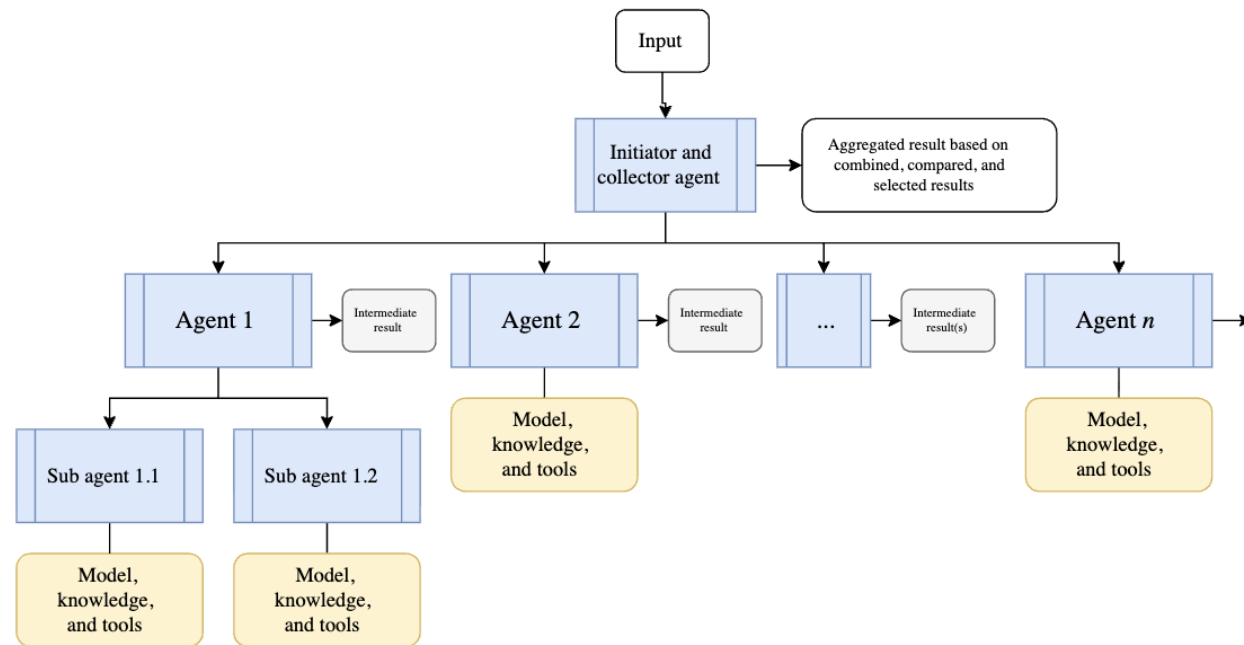
ABSTRACT

While large language models (LLMs) have demonstrated impressive performance across tasks in language understanding and interactive decision making, their abilities for reasoning (e.g. chain-of-thought prompting) and acting (e.g. action plan generation) have primarily been studied as separate topics. In this paper, we explore the use of LLMs to generate both reasoning traces and task-specific actions in an interleaved manner, allowing for greater synergy between the two: reasoning traces help the model induce, track, and update action plans as well as handle exceptions, while actions allow it to interface with and gather additional information from external sources such as knowledge bases or environments. We apply our approach, named ReAct, to a diverse set of language and decision making tasks and demonstrate its effectiveness over state-of-the-art baselines in addition to improved human interpretability and trustworthiness. Concretely, on question answering (HotpotQA) and fact verification (Fever), ReAct overcomes prevalent issues of hallucination and error propagation in chain-of-thought reasoning by interacting with a simple Wikipedia API, and generating human-like task-solving trajectories that are more interpretable than baselines without reasoning traces. Furthermore, on two interactive decision making benchmarks (ALFWorld and WebShop), ReAct outperforms imitation and reinforcement learning methods by an absolute success rate of 34% and 10% respectively, while being prompted with only one or two in-context examples.

# Multi-Agent Patterns

Much more flexible

- Large design space!
- Most common approach: central **orchestrator** splits up work and delegates to **specialists**

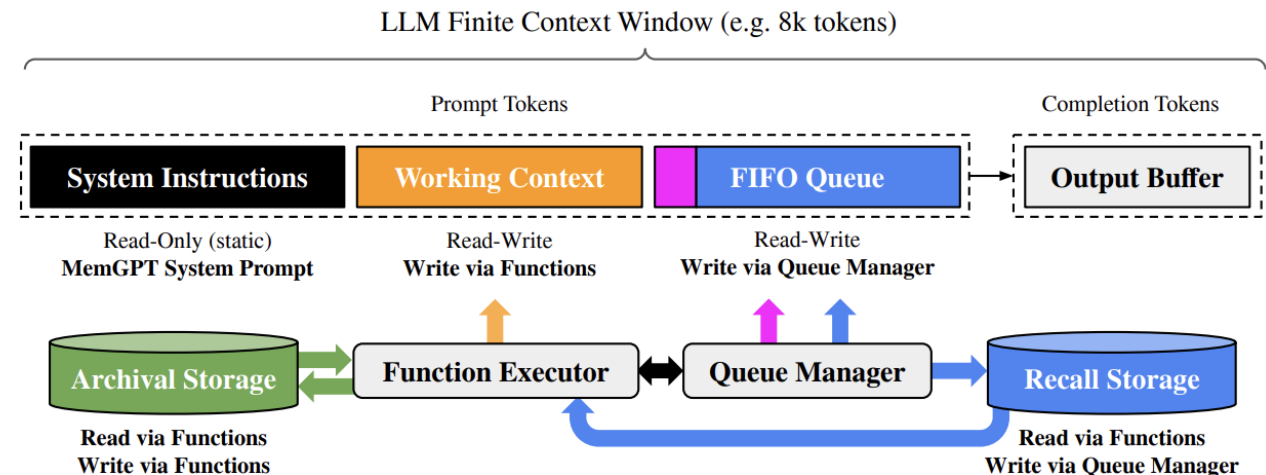


Microsoft

# Components: Memories

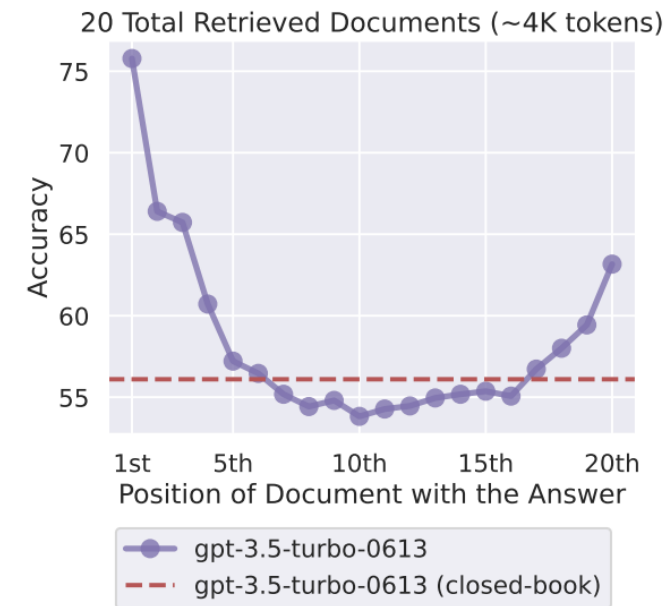
Interactive agents require the use of **memory**

- I.e., we need to store state/other useful information
- Simplest approach: **use LLM context window** as the memory
  - Increasingly tractable with very long-context windows
- Can also create memory systems/agents
  - E.g., MemGPT (Packer et al '24)
- Or buffers, RAG, etc.



# Interlude: Long Context Windows

- Simplest approach: **use LLM context window** as the memory
  - Increasingly tractable with very long-context windows
- Note that the advertised long context may not be realistic!
- Liu et al '23 (“Lost in the Middle: How Language Models Use Long Contexts”)
- Needle-in a haystack benchmarks
  - Numerous variants



# Components: Reasoning and Planning

Classic AI problem (decompose task, plan)

- Addressable via reasoning capabilities of LLMs,
  - Including older methods: chain-of-thought, tree-of-thoughts, self-reflection
- Example: **ReAct** (Reason + Act) (Yao et al '23)

(1) Hotspot QA

**Question:** Aside from the Apple Remote, what other device can control the program Apple Remote was originally designed to interact with?

(1d) ReAct (Reason + Act)

**Thought 1:** I need to search Apple Remote and find the program it was originally designed to interact with.

**Act 1:** `Search[Apple Remote]`

**Obs 1:** The Apple Remote is a remote control introduced in October 2005 by Apple ... originally designed to control the `Front Row media center` program ...

**Thought 2:** Apple Remote was originally designed to control the `Front Row media center` program. I need to search `Front Row` next and find what other device can control it.

**Act 2:** `Search[Front Row]`

**Obs 2:** Could not find [Front Row]. Similar: ['Front Row Seat to Earth', 'Front Row Motorsports', '`Front Row (software)`', ...]

**Thought 3:** `Front Row` is not found. I need to search `Front Row (software)` .

**Act 3:** `Search[Front Row (software)]`

**Obs 3:** `Front Row` is a discontinued media center software ...

**Thought 4:** `Front Row (software)` is controlled by an `Apple Remote` or the `keyboard function keys`. So the answer is `keyboard function keys`.

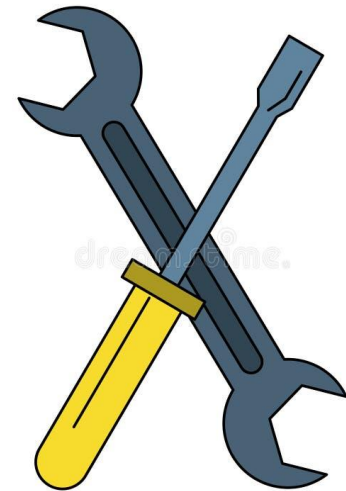
**Act 4:** `Finish[keyboard function keys]`



# Components: Tool Use

Standard component in agentic settings

- Where do we get tools?
  - APIs, web search, etc
  - Or, in some cases use of physical tools
- Most agent frameworks permit the use of web/code-based tools
- Major challenge: tool usage is **noisy**



# Tool Noise/Uncertainty

Lots of sources of error,

- Call the wrong tool, interpret output incorrectly, fail to chain together tools in the right way
- Ex: Agentic Insurance Underwriting Benchmark

**“Tool use errors:** Across models, including top performers, agents made at least one tool call error in 36% of the conversations despite access to the metadata required to use tools properly.”

```
===== tool_calls =====
----- tool -----
read_query

----- input -----
{"query": "SELECT DISTINCT NAICS_Code, Description FROM naics WHERE Description LIKE '%flour milling%' OR Description LIKE '%wheat flour%' OR Description LIKE '%corn flour%'"}

===== read_query =====
Error: ToolException('Error executing tool read_query: SQLite error: no such column: NAICS_Code') Please fix your mistakes.
```

# Creating New Tools

Note that LLMs can both use & **create** tools

- For example, LATM (Cai '24)

- Powerful model tool creator, less powerful is user

Tool maker side – only run *once* on a powerful model

## Tool proposing input

Input 1: Five animals are arranged in an unknown order.  
Conditions:  
🐶 > 🐱  
🐭 < 🐱  
🐱 is the second one  
🐭 > 🐹  
Question:  
What animal is the third one?  
Output 1: 🐭  
...  
Input k: \*\*\*  
Output k: \*\*\*  
  
Please write a generic Python function to solve this type of problems

## Tool proposing output

```
from itertools import permutations

def find_order(objects, conditions):
    for order in permutations(objects):
        valid = True
        for condition in conditions:
            if not condition(order):
                valid = False
                break
        if valid:
            return order
```

Tool user side – run *several times* on a lightweight model

## Tool using output

```
animals = ('🐶', '🐱', '🐭', '🐹', '🐭', '🐹')

condition1 = lambda order:
order.index('🐶') > order.index('🐱')
condition2 = lambda order:
order.index('🐭') < order.index('🐱')
condition3 = lambda order: order[1] == '🐭'
condition4 = lambda order:
order.index('🐭') > order.index('🐹')

conditions = (condition1, condition2,
condition3, condition4)

order = find_order(animals, conditions)

print(order[2])
```



**Break & Questions**

# Outline

- **Introduction to LLM-Powered Agents**
  - Motivation, goals, differences vs classical agents / standard LMs, overall architectures and key components
- **Multiagent Systems**
  - Motivation, simulations, architectures and design, communication
- **Challenges**
  - Scaling, realism, open problems

# Multi-Agent Systems

One advantage of agentic approach: can create multiple agents that cooperate to perform tasks. Why?

- **Scale** up workers, permit **specialization**
- Of course, inspired by earlier work,
  - Pre-LLM agent networks called “swarms”,
  - And distributed computing more broadly
- Also useful for simulations
  - Park et al ‘23



# Interlude: Societal Simulations

Popular new area,

- Difficult to perform large-scale studies of human behavior,
- Try to set up a simulation of how humans interact via agents

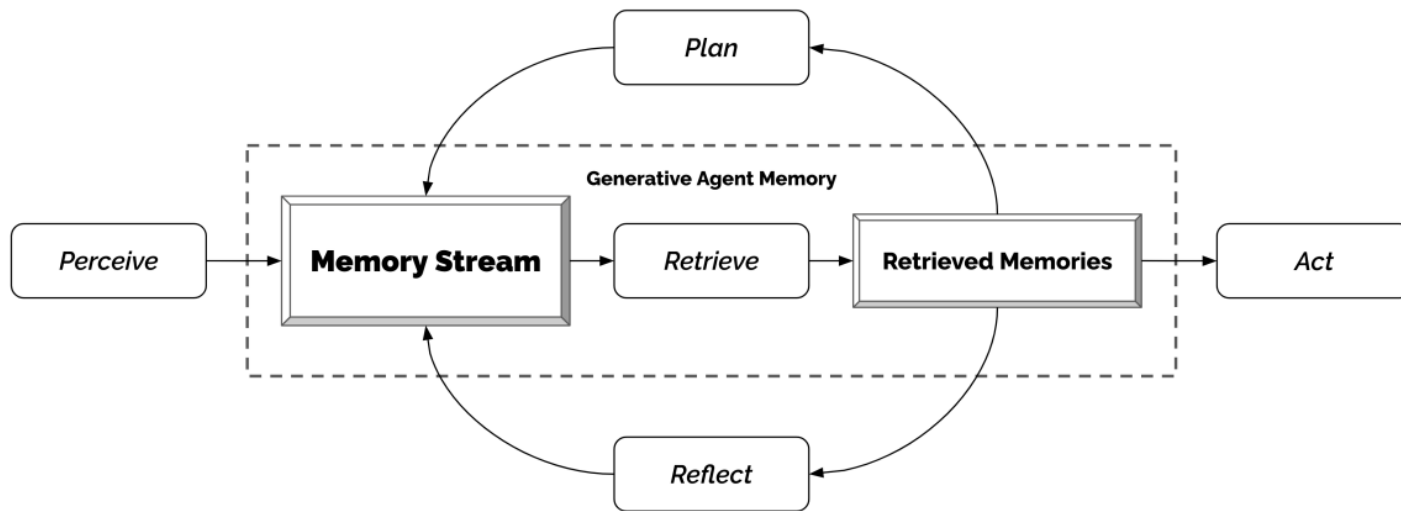


Park et al '23

# Interlude: Societal Simulations

Popular new area,

- Try to set up a simulation of how humans interact via agents
- Itself a large multi-agent system: must tackle each of the design questions



**Memory Stream**

```
2023-02-13 22:48:20: desk is idle
2023-02-13 22:48:20: bed is idle
2023-02-13 22:48:10: closet is idle
2023-02-13 22:48:10: refrigerator is idle
2023-02-13 22:48:10: Isabella Rodriguez is stretching
2023-02-13 22:33:30: shelf is idle
2023-02-13 22:33:30: desk is neat and organized
2023-02-13 22:33:10: Isabella Rodriguez is writing in her journal
2023-02-13 22:18:10: desk is idle
2023-02-13 22:18:10: Isabella Rodriguez is taking a break
2023-02-13 21:49:00: bed is idle
2023-02-13 21:48:50: Isabella Rodriguez is cleaning up the kitchen
2023-02-13 21:48:50: refrigerator is idle
2023-02-13 21:48:50: bed is being used
2023-02-13 21:48:10: shelf is idle
2023-02-13 21:48:10: Isabella Rodriguez is watching a movie
2023-02-13 21:19:10: shelf is organized and tidy
2023-02-13 21:18:10: desk is idle
2023-02-13 21:18:10: Isabella Rodriguez is reading a book
2023-02-13 21:03:40: bed is idle
2023-02-13 21:03:30: refrigerator is idle
2023-02-13 21:03:30: desk is in use with a laptop and some papers on it
...
```

**Q. What are you looking forward to the most right now?**

Isabella Rodriguez is excited to be planning a Valentine's Day party at Hobbs Cafe on February 14th from 5pm and is eager to invite everyone to attend the party.

retrieval	recency	importance	relevance
2.34	0.91	0.63	0.80

ordering decorations for the party

2.21	0.87	0.63	0.71
------	------	------	------

researching ideas for the party

2.20	0.85	0.73	0.62
------	------	------	------

...

I'm looking forward to the Valentine's Day party that I'm planning at Hobbs Cafe!

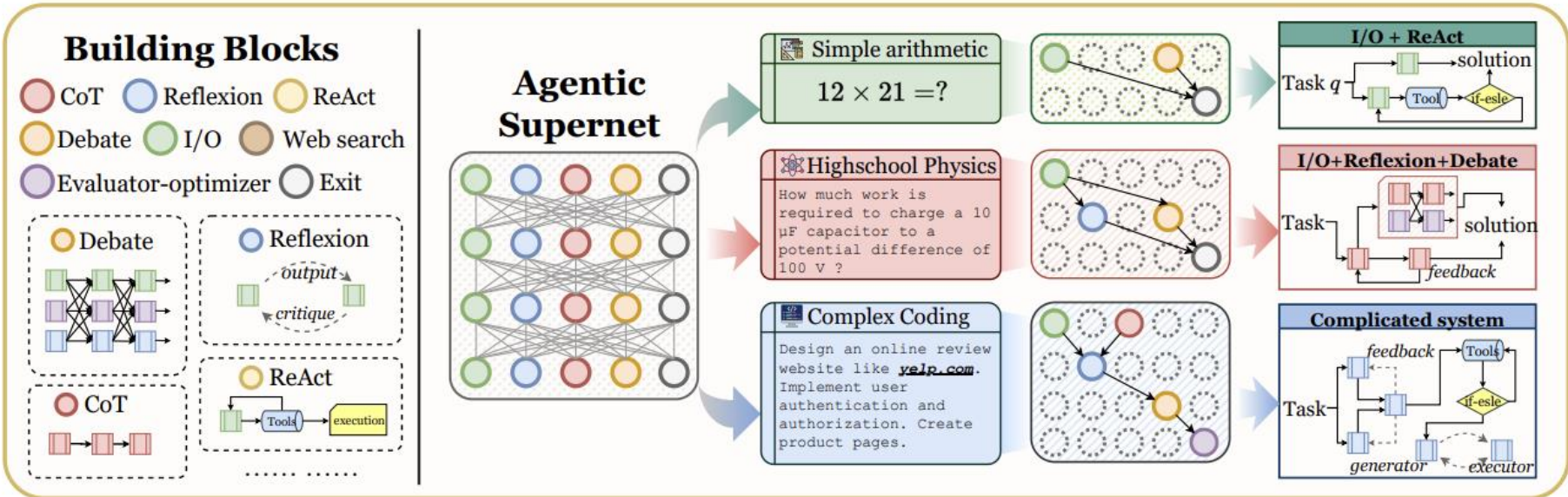


Isabella

# MAS Architecture Search

Lots of design choices to make!

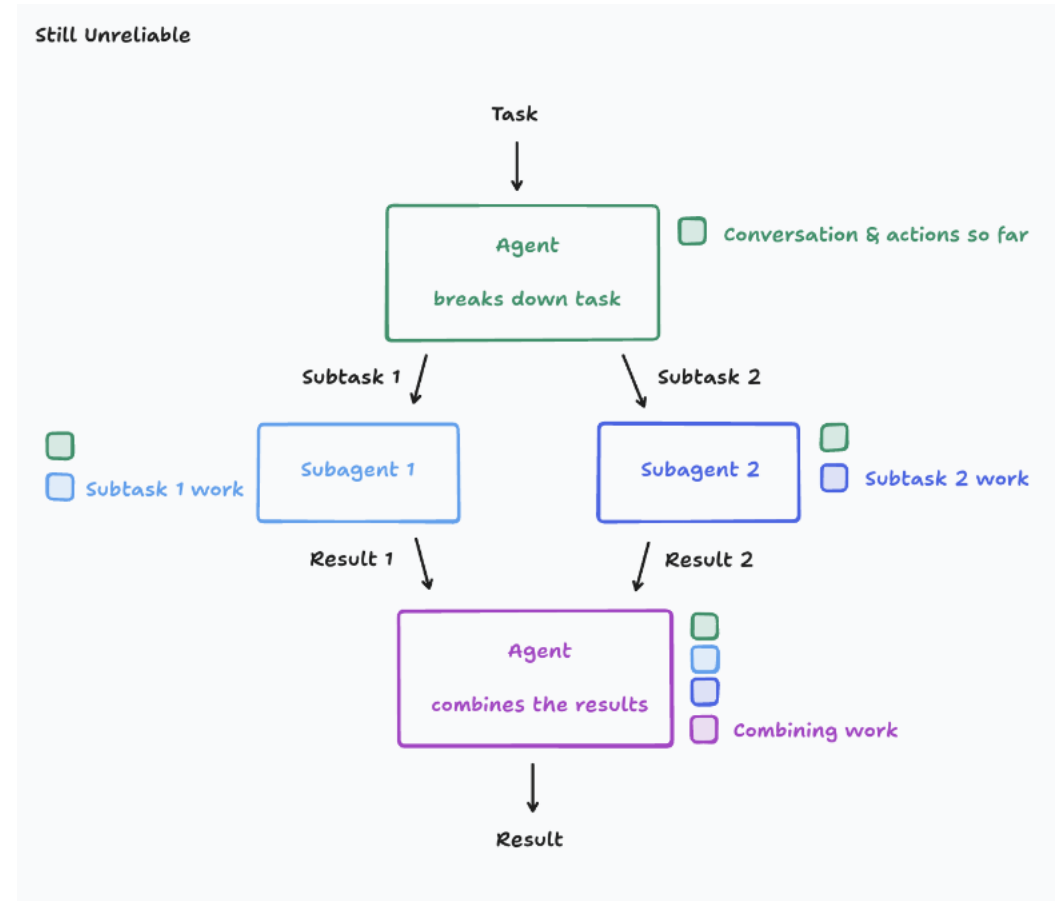
- Can borrow principles from **neural architecture search** to construct an overall agent architecture
- Example: Zhang et al '25, "Multi-agent Architecture Search via Agentic Supernet"



# Multi-Agent Systems: Communication

Critical: each agent has access to information that is consistent with others,

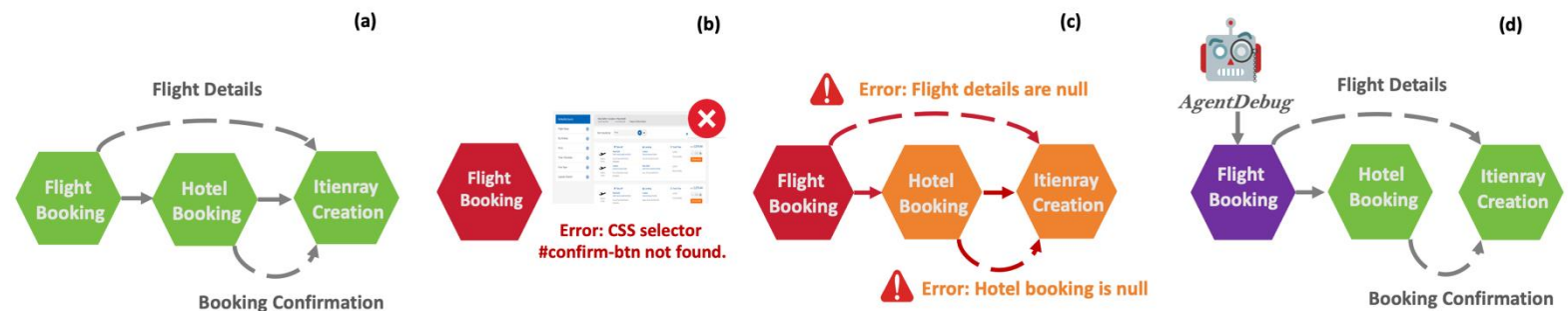
- But communicating everything between every pair of users is expensive.
- Nice blog on the challenges: **“Don’t Build Multi-Agents”** (Yan ’25 / Cognition)



# Multi-Agent Systems: Failure Modes

More agents → more can go wrong!

- Agents **repeat each other's mistakes**
- **Noisy** communication (and expensive!)
- Subgoals become **misaligned** with the main goal
- Central orchestrator aggregates weak outputs **incorrectly**
- Verification **bottlenecks**
  - Easier to take actions than verify them...





# Break & Questions

# Outline

- **Introduction to LLM-Powered Agents**
  - Motivation, goals, differences vs classical agents / standard LMs, overall architectures and key components
- **Multiagent Systems**
  - Motivation, simulations, architectures and design, communication
- **Challenges**
  - Scaling, realism, open problems

# Agentic Challenges

Things get more complex than just static LLMs...

- Large design and configuration space!
- Lots of places where things can go wrong too



# Failure Modes

Lots of possibilities! We saw some MAS-specific ones, but this is a more general problem with agents:

- **Hallucinated** tool calls
- **Memory drift** (summarization errors, staleness)
- **Brittleness** (small changes to UI or API change)
- **Misalignment** (optimize the wrong thing)
- **Cost and latency issues**

# Scaling Agents?

Kim et al '26: *Towards a Science of Scaling Agent Systems*

Tests out many configurations

- Vary benchmarks, agent archs., LLMs

Sample findings:

- **Benefits of MAS task dependent**
  - If not useful for tasks, burns tokens
  - Extra coordination can also hurt
- **Tool-heavy:** coordination harder
- **Central verification** useful

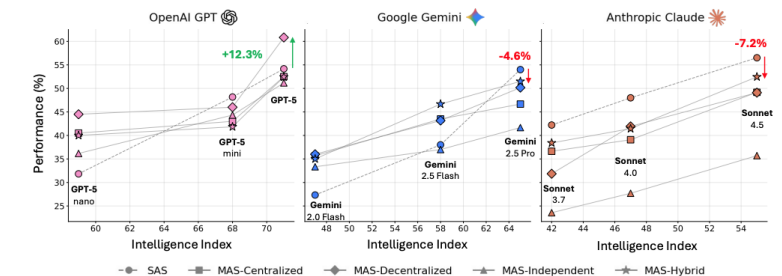
Google Research Google DeepMind

## Towards a Science of Scaling Agent Systems

Yubin Kim<sup>1,3,†</sup>, Ken Gu<sup>1</sup>, Chanwoo Park<sup>3</sup>, Chunjong Park<sup>2</sup>, Samuel Schmidgall<sup>2</sup>, A. Ali Heydari<sup>1</sup>, Yao Yan<sup>1</sup>, Zhihan Zhang<sup>1</sup>, Yuchen Zhuang<sup>2</sup>, Yun Liu<sup>1</sup>, Mark Malhotra<sup>1</sup>, Paul Pu Liang<sup>3</sup>, Hae Won Park<sup>3</sup>, Yuzhe Yang<sup>1</sup>, Xuhai Xu<sup>1</sup>, Yilun Du<sup>1</sup>, Shwetak Patel<sup>1</sup>, Tim Althoff<sup>1</sup>, Daniel McDuff<sup>1,†</sup> and Xin Liu<sup>1,†</sup>

<sup>1</sup>Google Research, <sup>2</sup>Google DeepMind, <sup>3</sup>Massachusetts Institute of Technology, <sup>†</sup>Corresponding Author

Agents, language model-based systems capable of reasoning, planning, and acting are widely adopted in real-world tasks, yet how their performance changes as these systems scale across key dimensions remains underexplored. We introduce quantitative *scaling principles* for agent systems as a predictive model, capturing how performance varies with coordination, model capability, and measurable system and task factors. Across 260 configurations spanning six agentic benchmarks, five canonical architectures (Single-Agent and four Multi-Agent: Independent, Centralized, Decentralized, Hybrid), and three LLM families, we perform controlled evaluations standardizing tools, prompts, and compute to isolate architectural effects. The resulting model achieves a cross-validated  $R^2=0.373$  across all six benchmarks ( $R^2=0.413$  with a task-grounded capability metric). We identify a robust capability-saturation effect and additional patterns: (1) a coordination yields diminishing returns once single-agent baselines exceed certain performance; (2) tool-heavy tasks appear to incur multi-agent overhead; and (3) architectures without centralized verification tend to propagate errors more than those with centralized coordination. Relative performance change compared to single-agent baseline ranges from +80.8% on decomposable financial reasoning to -70.0% on sequential planning, demonstrating that architecture-task alignment determines collaborative success. The framework identifies the best-performing architecture for 87% of held-out configurations and shows consistent relative architecture preferences on unseen frontier models. Agent effectiveness depends on alignment between coordination and task structure, and that mismatched coordination degrades the performance.



# Evaluating Memory

What about choice of memory?

- He et al '26 **MemoryArena**

- Explicitly tests performance in memory-bound tasks

- Benchmarks various choices:

- Long-context buffer,
- RAG system,
- Memory agent.

Findings: memory-bound tasks **still hard!**

---

## Benchmarking Agent Memory in Interdependent Multi-Session Agentic Tasks

---

Zexue He<sup>\*1</sup> Yu Wang<sup>\*2</sup> Churan Zhi<sup>\*2</sup> Yuanzhe Hu<sup>\*2</sup> Tzu-Ping Chen<sup>\*2</sup> Lang Yin<sup>\*3</sup> Ze Chen<sup>4</sup>  
Tong Arthur Wu<sup>5</sup> Siru Ouyang<sup>3</sup> Zihan Wang<sup>6</sup> Jiaxin Pei<sup>1</sup> Julian McAuley<sup>2</sup> Yejin Choi<sup>1</sup> Alex Pentland<sup>1</sup>

### Abstract

Existing evaluations of agents with memory typically assess *memorization* and *action* in isolation. One class of benchmarks evaluates memorization by testing recall of past conversations or text but fails to capture how memory is used to guide future decisions. Another class focuses on agent acting in single-session tasks without the need for long-term memory. However, in realistic settings, memorization and action are tightly coupled: agents acquire memory while interacting with the environment, and subsequently rely on that memory to solve future tasks. To capture this setting, we introduce MEMORYARENA, a unified evaluation gym for benchmarking agent memory in *multi-session Memory-Agent-Environment loops*. The benchmark consists of human-crafted agentic tasks with explicitly interdependent subtasks, where agents must learn from earlier actions and feedback by distilling experiences into memory, and subsequently use that memory to guide later actions to solve the overall task. MEMORYARENA supports evaluation across web navigation, preference-constrained planning, progressive information searching, and sequential formal reasoning, and reveals that agents with near-saturated performance on existing long-context memory benchmarks like LoCoMo perform poorly in our agentic setting, exposing a gap in current evaluations for agents with memory. MEMORYARENA is released at <https://memoryarena.github.io/>.

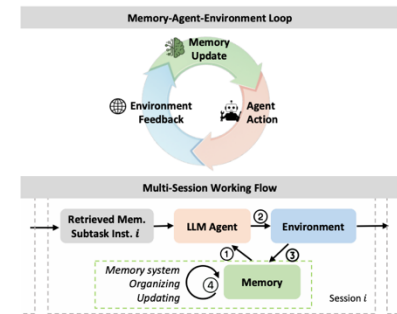


Figure 1. MEMORYARENA Evaluates agents with Memory with multi-session tasks in a Memory-Agent-Environment Loop.

through interaction with an environment (*action*) (Hu et al., 2025b). However, existing evaluations of LLM agents with memory typically isolate and assess only one aspect. The first class of benchmarks focuses on evaluating *memorization* through recall or retrieval over static long-context inputs in question answering or summarization settings (Wu et al., 2025; Zhong et al., 2024; Maharana et al., 2024; Hu et al., 2025b), including benchmarks such as LoCoMo (Maharana et al., 2024) and LongMemEval (Wu et al., 2025). In these setups, agents are required to memorize provided conversations or text chunks, and are evaluated on whether they can recall specific information through downstream QA tasks. However, despite being effective at measuring factual recall, such benchmarks do not involve agentic decision-making, environment dynamics, or action-dependent con-



**Thank You!**