

Developing Agentic Applications

Workshop: Multi-Agent Applications



Henry Ruiz, PhD.
Research Scientist
Texas A&M University
GDE AI & GCP
@devharuiz
<https://haruiz.github.io/>



Juan Guillermo, M.Sc.
Tech Lead @ WordBox
GDE AI & GCP & Firebase
@jggomezt
<https://devhack.co>

Agenda

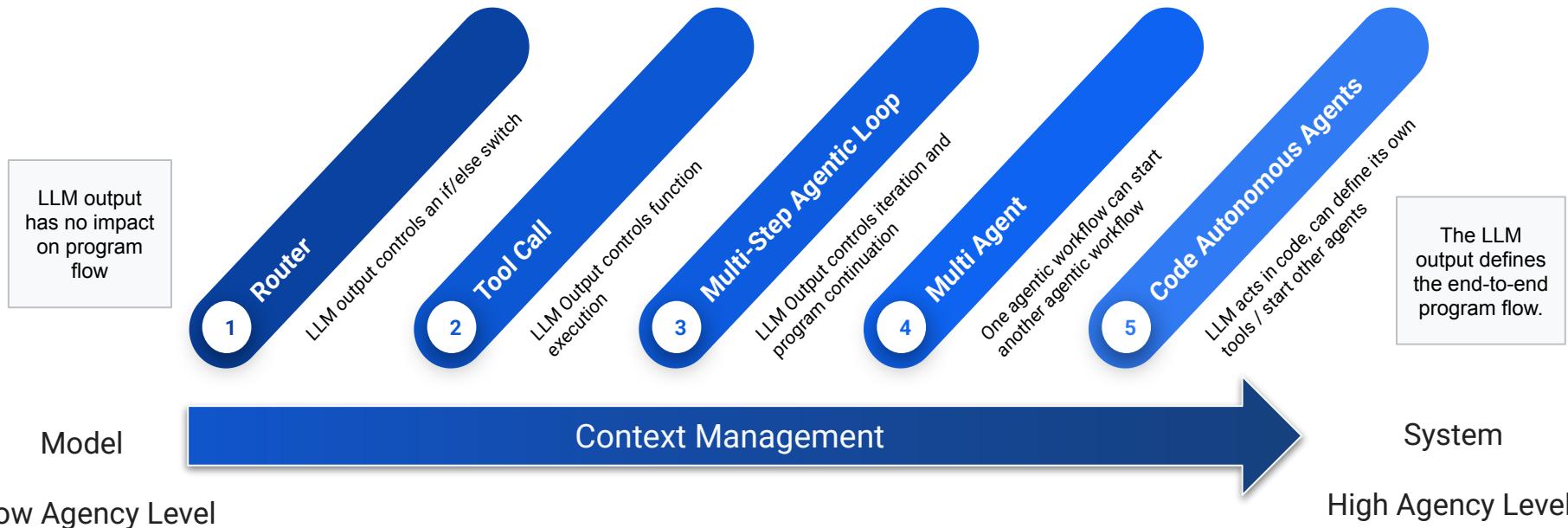
- Recap from previous section
- Emerging Capabilities of LLMs
(Reasoning, Planning, etc.)
- Controlling and Aligning LLMs
(State-of-the-Art Post-Pre Training
Techniques)
- From Chat Applications to Agents
- Anatomy of an agentic application
- Frameworks for Agent Development
- Multi Agent Applications Patterns

GCP Credits

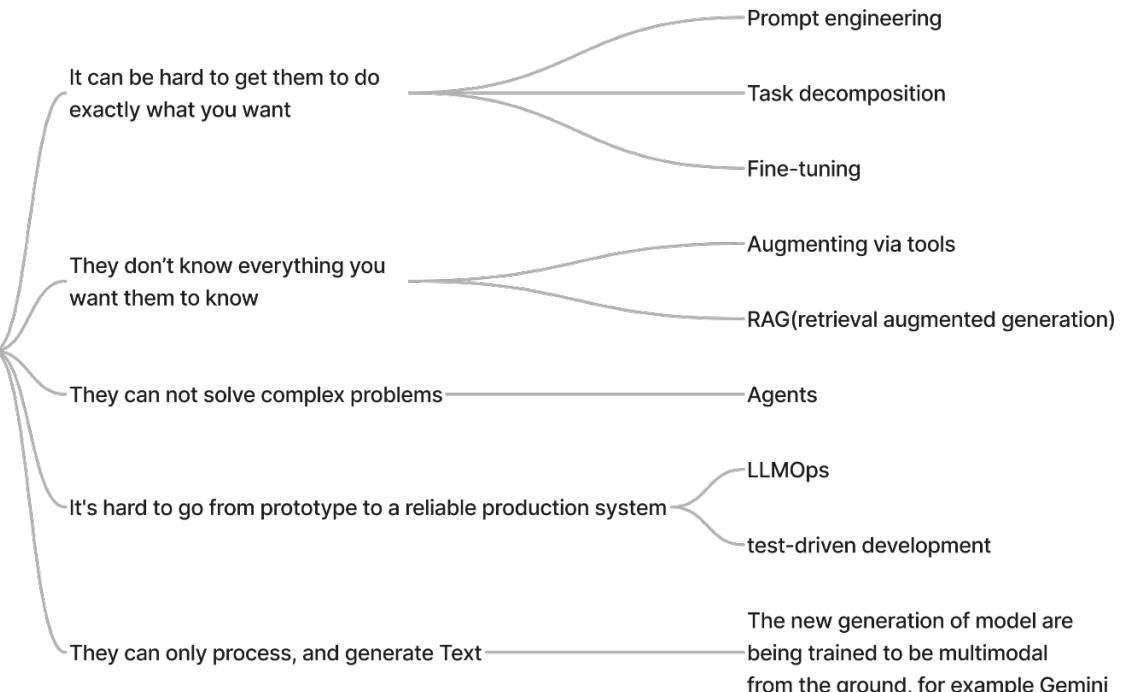


From our Previous section....

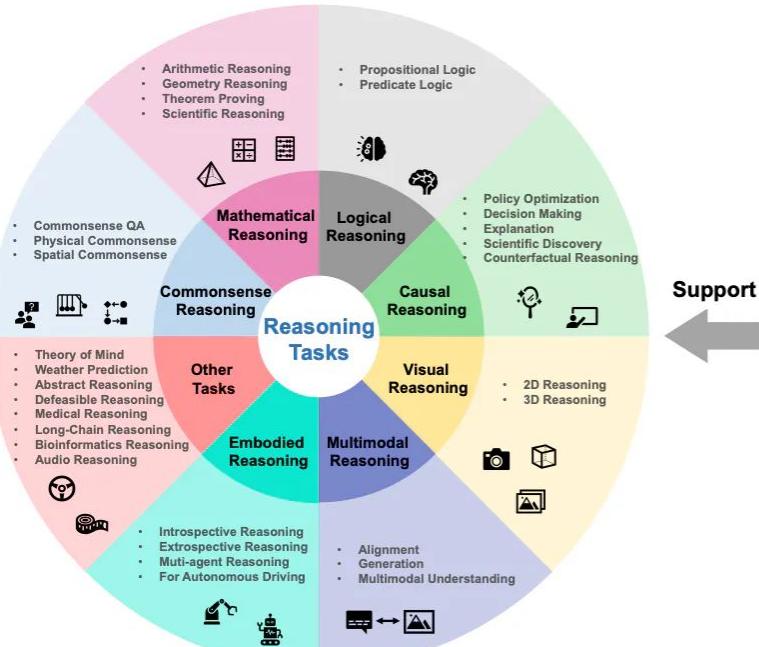
“agent” is not a discrete, 0 or 1 definition: instead, “agency” evolves on a continuous spectrum, as you give more or less power to the LLM on your workflow.



LLMs are magical, but they hallucinate (knowledge compression)



Recent advances in Large Language Models (LLMs) have demonstrated impressive reasoning capabilities. LLMs can convert unstructured text into structured outputs, enabling the automation of a wide range of tasks across diverse domains



Dictionary
Definitions from Oxford Languages · [Learn more](#)

re·a·son·ing
/rēzənİNG, 'rēznlNG/

noun

the action of thinking about something in a logical, sensible way.
"he explained the **reasoning** behind his decision at a media conference"

APA Dictionary of Psychology

Search and select a Dictionary term

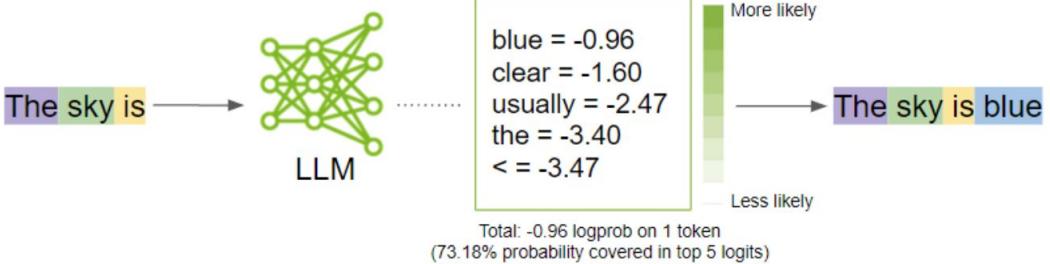
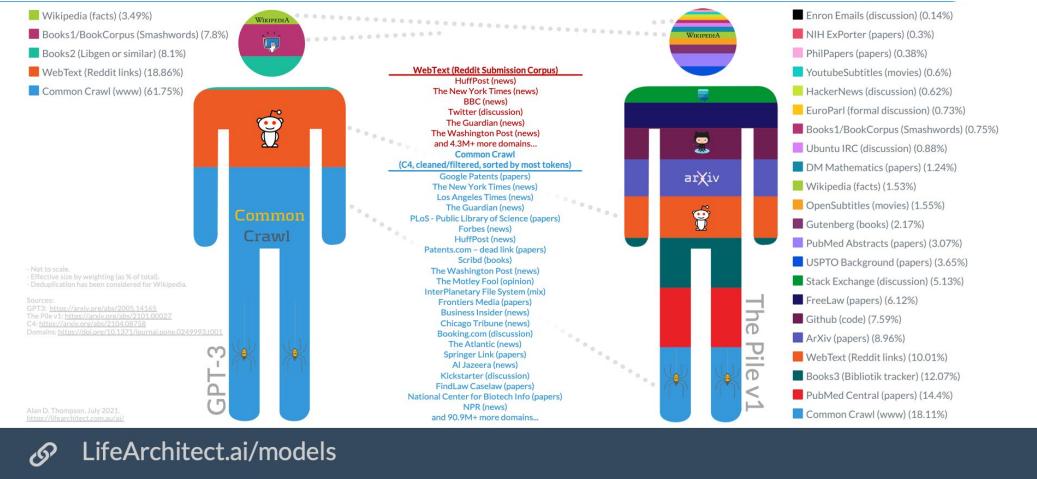
reasoning

Updated on 04/19/2018

n.

- thinking in which logical processes of an inductive or deductive character are used to draw conclusions from facts or premises. See **deductive reasoning**; **inductive reasoning**.
- the sequence of arguments or proofs used to establish a conclusion in this way. —**reason** vb.

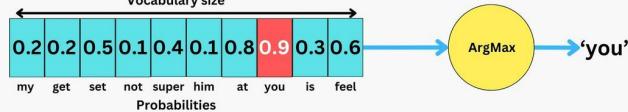
How Do Large Language Models (LLMs) Work?



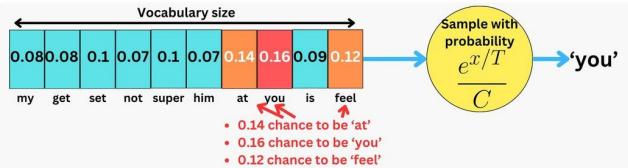
How LLMs Generate Text

TheAiEdge.io

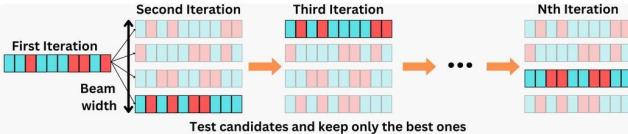
The Greedy strategy: Choose the token with highest probability



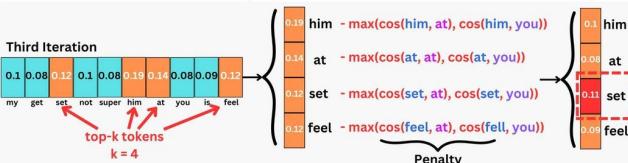
The Multinomial sampling strategy: Sampling tokens by using the probability



Beam Search: Maximize probability of the whole sequence

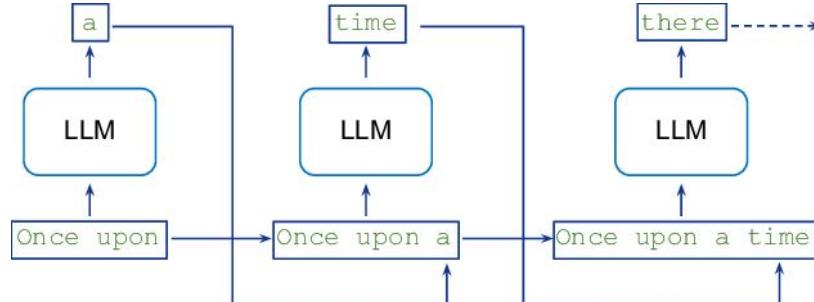


Contrastive search: Penalize repetitiveness



Credits to: <https://newsletter.theaiedge.io/about>

The autoregressive nature of LLMs suggests that they are capable of reasoning



Autoregressive (AR) language models:

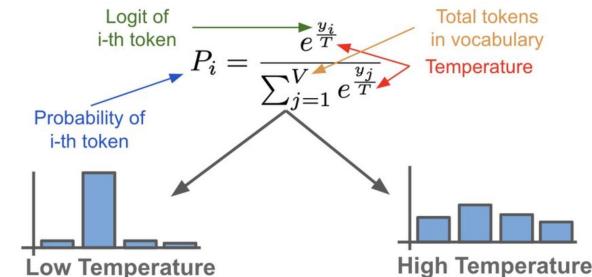
$$p(x_1, \dots, x_L) = p(x_1)p(x_2|x_1)p(x_3|x_2, x_1)\dots = \prod_i p(x_i|x_{1:i-1})$$

No approx: chain rule of probability

=> You only need a model that can predict the next token given past context!

Stanford

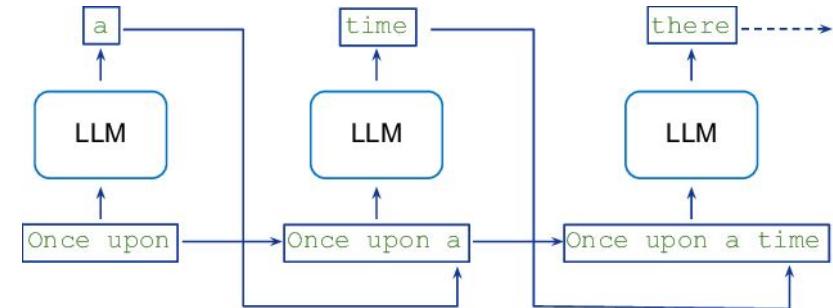
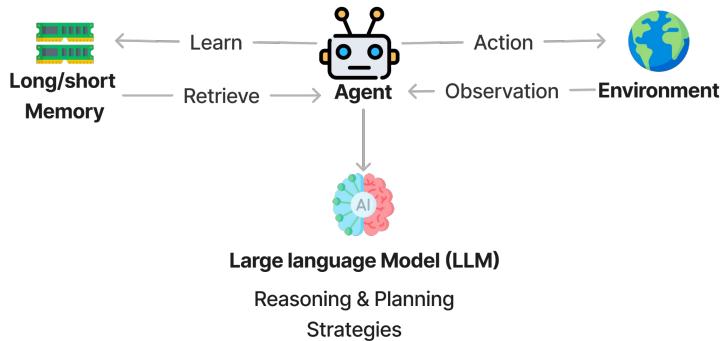
We can adjust the temperature to modulate the uniformity of the token distribution produced by the softmax transformation



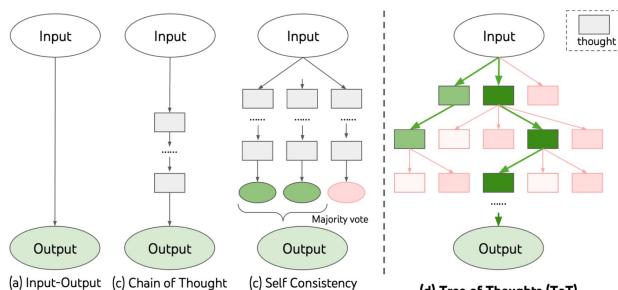
$$\arg \max \mathbb{P}(\text{final answer} | \text{problem})$$

$$= \sum_{\text{reasoning path}} \mathbb{P}(\text{ reasoning path, final answer} | \text{problem})$$

Key idea of reasoning: Derive the final answer through intermediate steps & Post training alignment



$$\begin{aligned} & \arg \max \mathbb{P}(\text{final answer} | \text{problem}) \\ = & \sum_{\text{reasoning path}} \mathbb{P}(\text{reasoning path}, \text{final answer} | \text{problem}) \end{aligned}$$

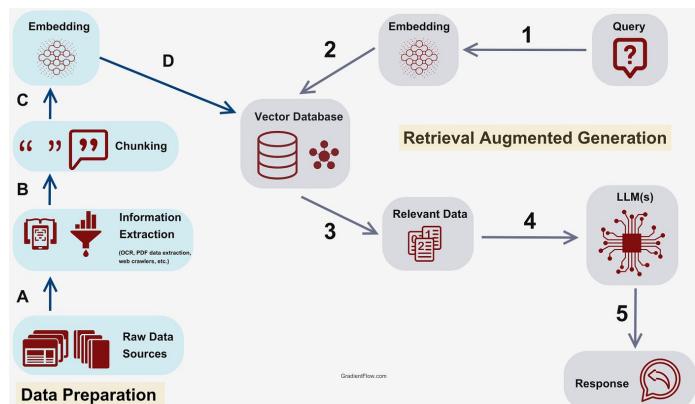


Wei Jason et al. "ReAct: Synergizing Reasoning and Acting in Large Language Models." *ArXiv:2201.11903 [Cs]*, 10 Oct. 2022, arxiv.org/abs/2201.11903

Yao, Shunyu, et al. "ReAct: Synergizing Reasoning and Acting in Language Models." *ArXiv.org*, 9 Mar. 2023, arxiv.org/abs/2210.03629.

Ling, Wang, et al. "Program Induction by Rationale Generation : Learning to Solve and Explain Algebraic Word Problems." *ArXiv:1705.04146 [Cs]*, 23 Oct. 2017, arxiv.org/abs/1705.04146.

Today, we can only certainly say that LLMs are universal approximators for retrieval tasks. This capability is supported by their contextual understanding and the autoregressive des



<https://arxiv.org/pdf/2310.07554>

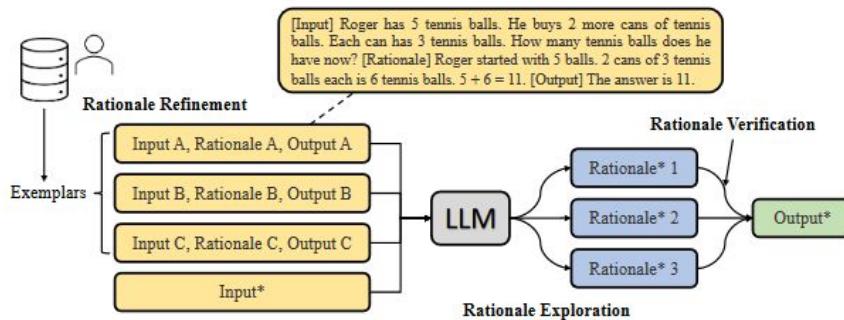
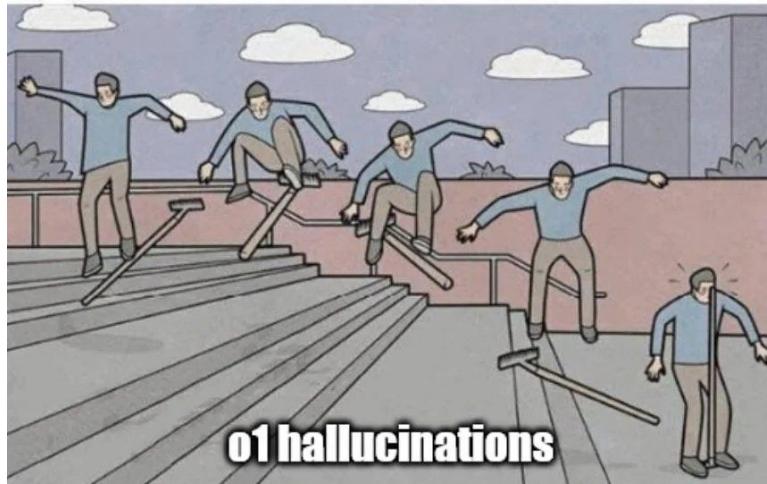


Figure 2: An illustration of *Chain-of-Thought Prompting and Rationale Engineering*, where asterisk (*) denotes the target problem to be solved.

<https://arxiv.org/pdf/2403.04121>

**GPT-4o
hallucinations**



o1 hallucinations

These capabilities has revolutionized the way we interact with and utilize Large Language Models (LLMs) today.

Q: What is $1 + 2$?



A: 3

Q: Which crop, cotton or soybeans, is better suited for the hot and dry climate of West Texas?



Requires knowledge

Q: A farmer has 5 hectares of land and wants to apply 150 kg of fertilizer per hectare. How much total fertilizer does the farmer need?



Requires Computation

Q: A farmer in Texas has two fields: Field A receives 30% more rainfall than Field B. If Field B receives 500 mm of rainfall per year, how much rainfall does Field A receive?



Requires reasoning

Retrieval augmented generation for knowledge (RAG)

Code augmentation

Tool use
(Nested calls)

What if both knowledge and reasoning are needed?

However, it is still a controversial claim

Can Large Language Models Reason and Plan?

Towards Reasoning in Large Language Models: A Survey

Jie Huang Kevin Chen-Chuan Chang

Department of Computer Science, University of Illinois at Urbana-Champaign
{jeffhj, kcchang}@illinois.edu

However, despite the strong performance of LLMs on certain reasoning tasks, it remains unclear whether LLMs are actually reasoning and to what extent they are capable of reasoning. For example, Kojima et al. (2022) claim that “LLMs are decent zero-shot reasoners (p. 1)”, while Valmeekam et al. (2022) conclude that “LLMs are still far from achieving acceptable performance on common planning/reasoning tasks which pose no issues for humans to do (p. 2).” This limitation is also stated by Wei et al. (2022b):

<https://arxiv.org/pdf/2403.04121>

<https://arxiv.org/abs/2212.10403>

Subbarao Kambhampati

School of Computing & Augmented Intelligence Arizona State University email: rao@asu.edu

(A version appears in the Annals of The New York Academy of Sciences:

<https://nyaspubs.onlinelibrary.wiley.com/doi/10.1111/nyas.15125>

To summarize, nothing that I have read, verified, or done gives me any compelling reason to believe that LLMs do reasoning/planning, as normally understood. What they do instead, armed with web-scale training, is a form of universal approximate retrieval, which, as I have argued, can sometimes be mistaken for reasoning capabilities. LLMs do excel in idea generation for any task-including those involving reasoning, and as I pointed out, this can be effectively leveraged to support reasoning/planning in LLM-Modulo frameworks⁶. In other words, LLMs already have enough amazing approximate retrieval abilities that can be gainfully leveraged, that we don't need to ascribe questionable reasoning/planning capabilities to them.^{††}

Study with infants suggests language not necessary for reasoning ability

by Bob Yirka , Medical Xpress

A team of researchers from Spain, Hungary and Poland has found via a study with infants that language may not be a necessity for the ability to reason. In their paper published in the journal *Science*, the group describes their study and what they found, and also offer some opinions on how their findings might be used to better understand the ability to reason. Justin Halberda with Johns Hopkins University offers a Perspective piece on the work done by the team in the same journal issue.

The Illusion of Thinking: Understanding the Strengths and Limitations of Reasoning Models via the Lens of Problem Complexity

Parshin Shojaee^{*†} Iman Mirzadeh^{*} Keivan Alizadeh
Maxwell Horton Samy Bengio Mehrdad Farajtabar



Abstract

Recent generations of frontier language models have introduced Large Reasoning Models (LRMs) that generate detailed thinking processes before providing answers. While these models demonstrate improved performance on reasoning benchmarks, their fundamental capabilities, scaling properties, and limitations remain insufficiently understood. Current evaluations primarily focus on established mathematical and coding benchmarks, emphasizing final answer accuracy. However, this evaluation paradigm often suffers from data contamination and does not provide insights

A comment on Shojaee et al. (2025)

The Illusion of the Illusion of Thinking

A Comment on Shojaee et al. (2025)

C. Opus* A. Lawsen†

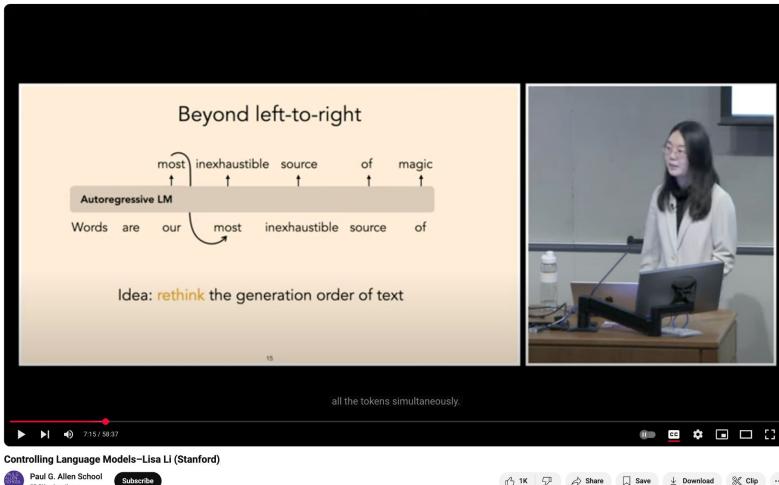
June 10, 2025

Abstract

Shojaee et al. (2025) report that Large Reasoning Models (LRMs) exhibit "accuracy collapse" on planning puzzles beyond certain complexity thresholds. We demonstrate that their findings primarily reflect experimental design limitations rather than fundamental reasoning failures. Our analysis reveals three critical issues: (1) Tower of Hanoi experiments systematically exceed model output token limits at reported failure points, with models explicitly acknowledging these constraints in their outputs; (2) The authors' automated evaluation framework fails to distinguish between reasoning failures and practical constraints, leading to misclassification of model capabilities; (3) Most concerningly, their River Crossing benchmarks include mathematically impossible instances for $N \geq 6$ due to insufficient boat capacity, yet models are scored as failures for not solving these unsolvable problems. When we control for these experimental artifacts, by requesting generating functions instead of exhaustive move lists, preliminary experiments across multiple models indicate high accuracy on Tower of Hanoi instances previously reported as complete failures. These findings highlight the importance of careful experimental

LLM Post-Training

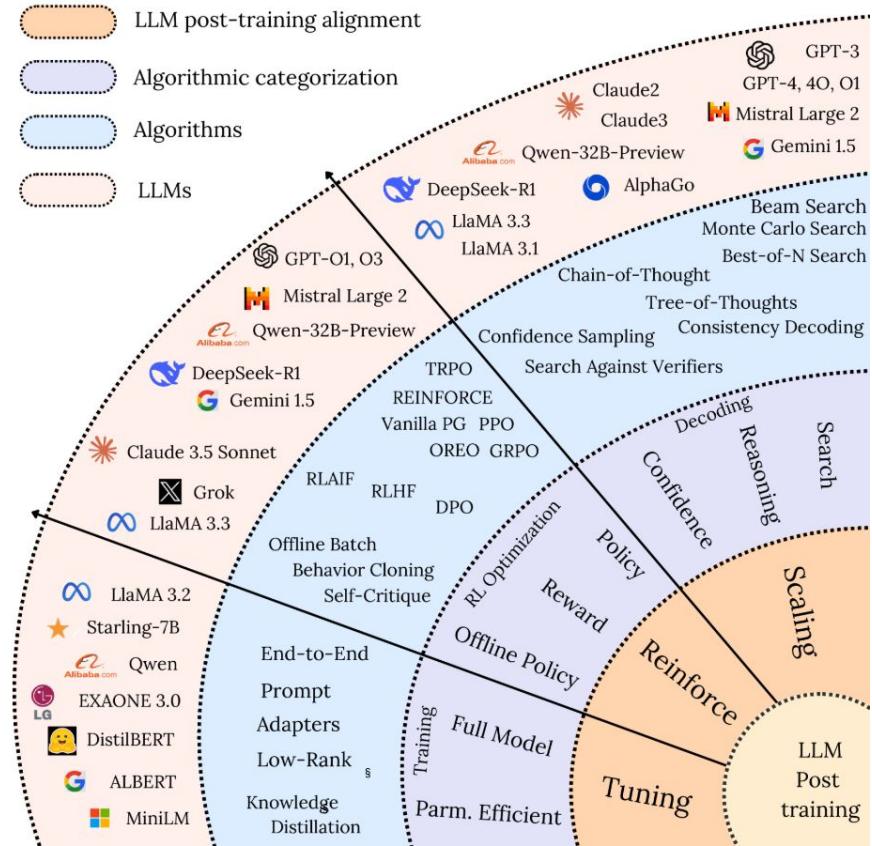
<https://www.youtube.com/watch?v=tEQ9N5JjGW0&t=3204s>



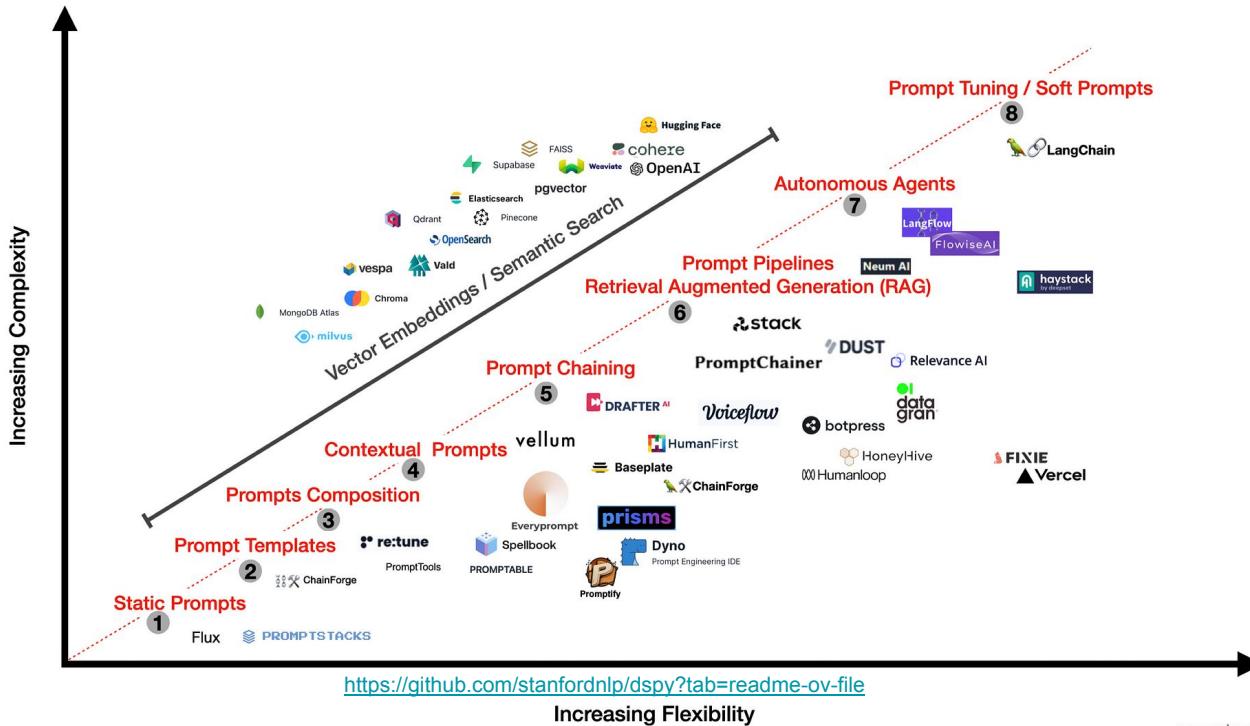
Controlling Language Models - Lisa Li (Stanford)

Paul G. Allen School
23.3K subscribers

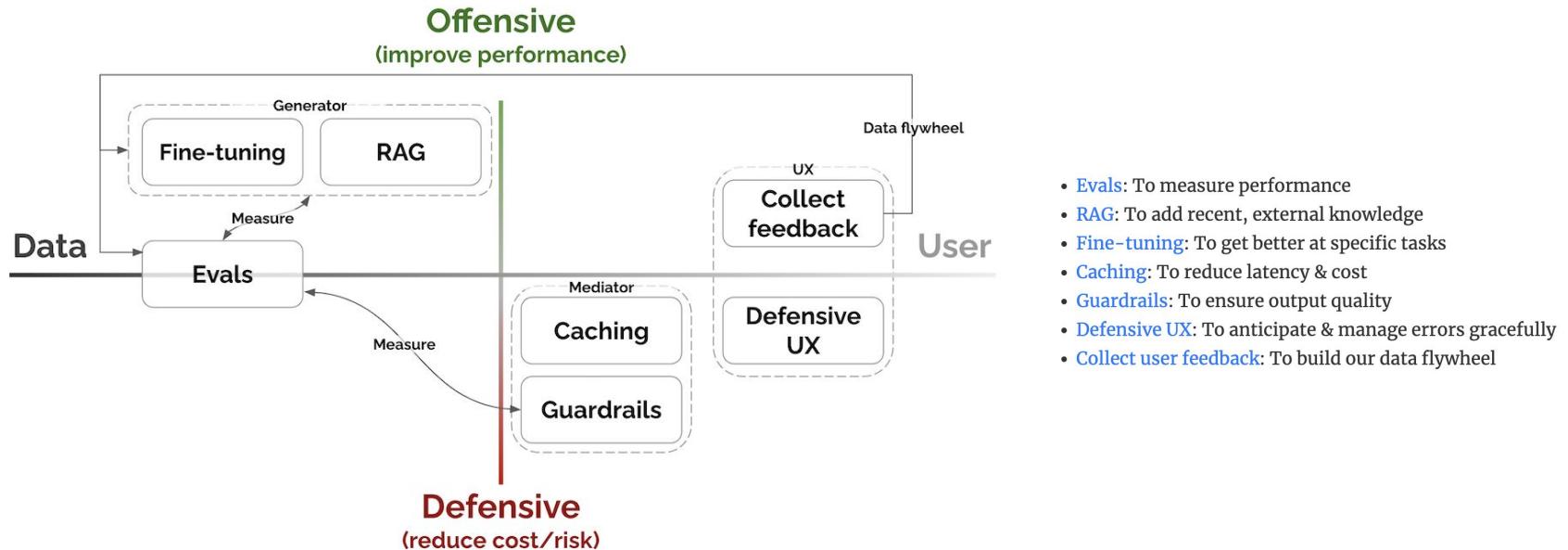
1K Share Save Download Clip ...



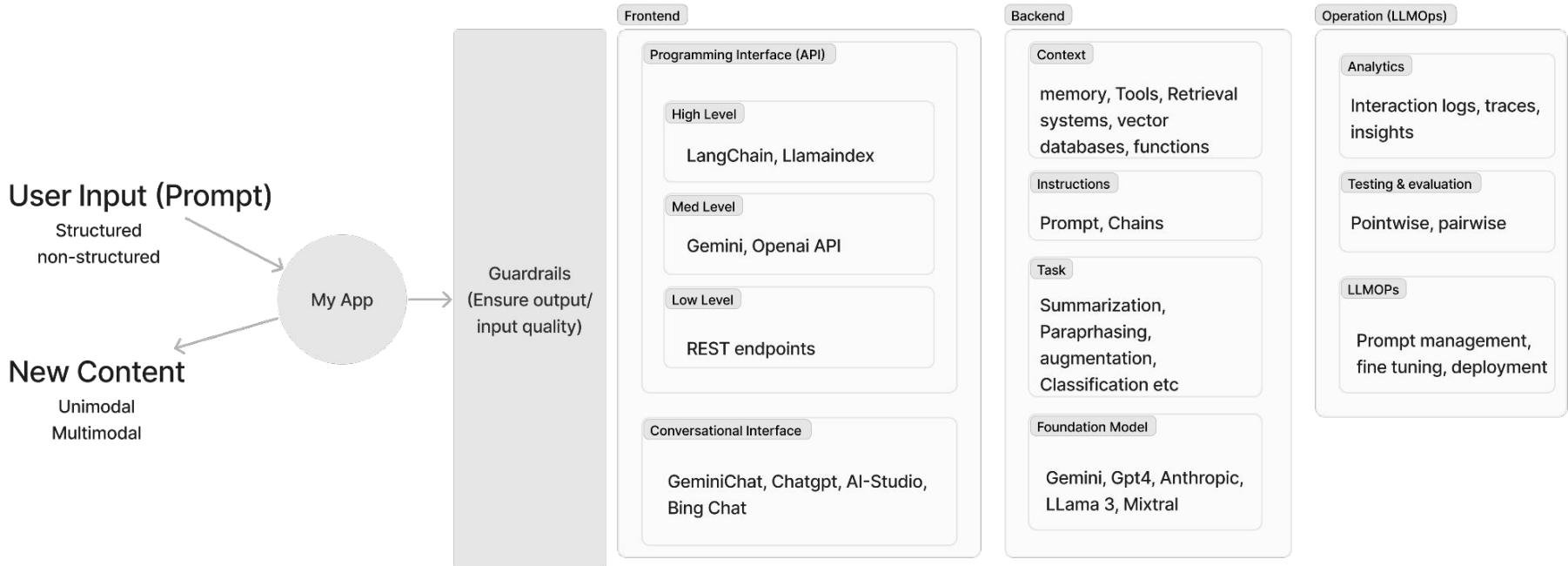
The challenges have opened the door to a new wave of emerging architectures and design patterns for developing LLM applications.



The GOAL: Improving LLMs performance



Anatomy of an LLM app



The key: Memory = Context Window

Retrieve Anything To Augment Large Language Models

<https://arxiv.org/pdf/2310.07554>

Augmented language models

...In the following slides, we'll explore how taking full advantage of the expanded context window in LLMs unlocks capabilities that go far beyond simple memorization—enabling more sophisticated reasoning, task automation, and contextual understanding

Retrieval



Augment with
a bigger corpus

Chains



Augment with
more LLM calls

Tools



Augment with
outside sources

Context provides LLMs with unique, up-to-date information.

...But it only fits a limited amount of information

← r/OpenAI · 6 mo. ago
veleros

... al

Gemini's context window is much larger than anyone else's

Image

Google's AI Gemini can fit almost 10 Harry Potter books in its context window

Number of copies of 'Harry Potter and the Sorcerer's Stone' books that can fit within each AI's context window

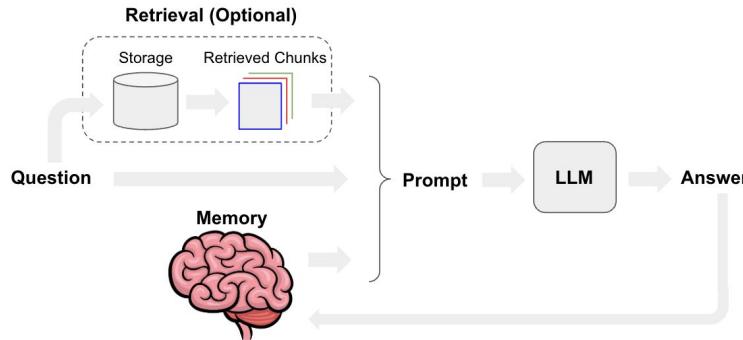
Model	Context Window Capacity
Gemini 1.5 Pro	9.77
Claude-2.1	1.95
gpt-4-turbo	1.25
Grok	0.24
gpt-3.5-turbo	0.16

The chart illustrates the context window capacity of various AI models by showing the number of 'Harry Potter and the Sorcerer's Stone' books that can fit within their respective context windows. Gemini 1.5 Pro has the largest capacity at 9.77, followed by Claude-2.1 at 1.95, gpt-4-turbo at 1.25, Grok at 0.24, and gpt-3.5-turbo with the smallest capacity at 0.16.

linkedin.com/in/eduardoviteri

Chat application

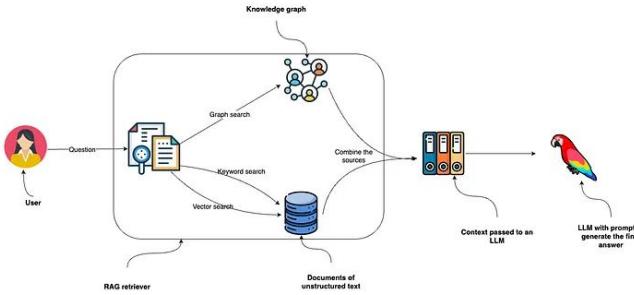
- The chat history is injected into the prompt.
- This provides the LLM (Large Language Model) with access to previous messages, referred to as "memory."
- The LLM uses these messages as context to generate its output.
- Alternative we can connect a chat app with a retrieval system



```
1 def get_chat_history():
2     """
3         This function returns the chat history.
4         :param history:
5         :return:
6         """
7     return [
8         HumanMessage("What is the capital of France?"),
9         AIMessage("The capital of France is Paris."),
10    ]
11
12 template = (
13     "Given the following conversation:"
14     "\nChat History:\n"
15     "{history}"
16     "\n, and a follow-up question:\n"
17     "{question}\n"
18     "respond to the follow-up question in english:"
19 )
20 prompt = ChatPromptTemplate.from_template(template)
21 llm = ChatGoogleGenerativeAI(model="gemini-pro")
22 chat_history = get_chat_history()
23
24 chat_context = RunnablePassthrough.assign(history=lambda _: chat_history)
25 chain = chat_context | prompt | llm | {"answer": StrOutputParser()}
26
27 response = chain.invoke({"question": "Can you respond again?"})
28 print(response)
```

RAG application (Vector index/ knowledge graph)

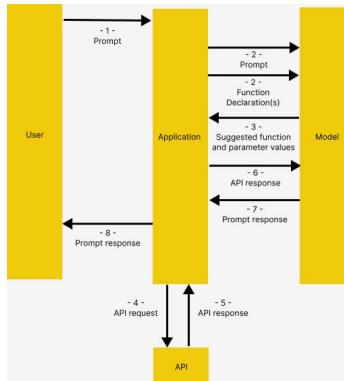
- Documents are transformed into embeddings.
- These embeddings are stored in an embedding database to create a vector index or a knowledge graph.
- When querying:
 - Relevant documents are retrieved from the vector index based on the prompt.
 - The LLM uses these documents as context to generate its output.



```
1 def get_rag_retriever():
2     """
3         This function returns the context for the RAG model.
4     :return:
5     """
6     loader = TextLoader("state_of_the_union.txt")
7     documents = loader.load()
8     text_splitter = CharacterTextSplitter(chunk_size=1000, chunk_overlap=0)
9     docs = text_splitter.split_documents(documents)
10    embeddings_model = GoogleGenerativeAIEMBEDDINGS(
11        model="models/embedding-001", task_type="retrieval_document"
12    )
13    db = FAISS.from_documents(docs, embeddings_model)
14    print(db.index.ntotal)
15    return db.as_retriever()
16
17
18 template = """Answer the question based only on the following context:
19 {context}
20 Question: {question}
21 """
22 prompt = ChatPromptTemplate.from_template(template)
23 llm = GoogleGenerativeAI(model="gemini-pro")
24 retriever = get_rag_retriever()
25
26 chat_context = RunnablePassthrough.assign(context=itemgetter("question") | retriever)
27 chain = chat_context | prompt | llm | {"answer": StrOutputParser()}
28 response = chain.invoke({"question": "what the state of the union talk about?"})
29 # response = chain.invoke({ "question": "talk me about apples?"})
30 print(response)
```

Function Calling

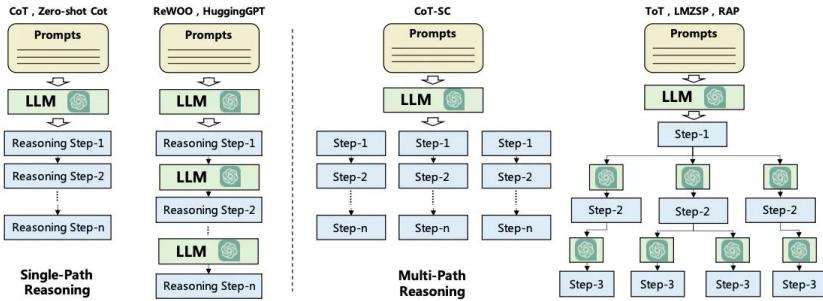
- A list of functions defined:
 - The signature of each function, including its description and arguments, is injected into the prompt.
 - This context provides access to a list of functions for delegating data processing tasks.
- The model does not directly call the functions.
 - Instead, it provides a structured data output.
 - This output includes the name of a selected function and the arguments that the model proposes for the function to be called with.



```
1 """
2 Basic example of using function calling with Gemini
3 :return:
4 """
5
6     # Create a model
7 def multiply(a: float, b: float):
8     """returns a * b."""
9     return a * b
10
11 def divide(a: float, b: float):
12     """returns a / b."""
13     return a / b
14
15 model = genai.GenerativeModel(
16     model_name="gemini-pro",
17     tools=[multiply,divide]
18 )
19 chat = model.start_chat(enable_automatic_function_calling=True)
20 response = chat.send_message(
21     "I have 57 cats, each owns 44 mittens, how many mittens is that in total?"
22 )
23 console.print(response.text)
24 for content in chat.history:
25     part = content.parts[0]
26     print(content.role, ">", type(part).to_dict(part))
27     print("-" * 80)
```

Agents

- LLMs of today (e.g., Gemini, GPT4, LLAMA 3, etc.):
- When given access to tools (e.g., external memory, a calculator, web search, code execution, etc.),
 - And context, Become agents that can **reason** and **act** on our behalf to solve tasks.

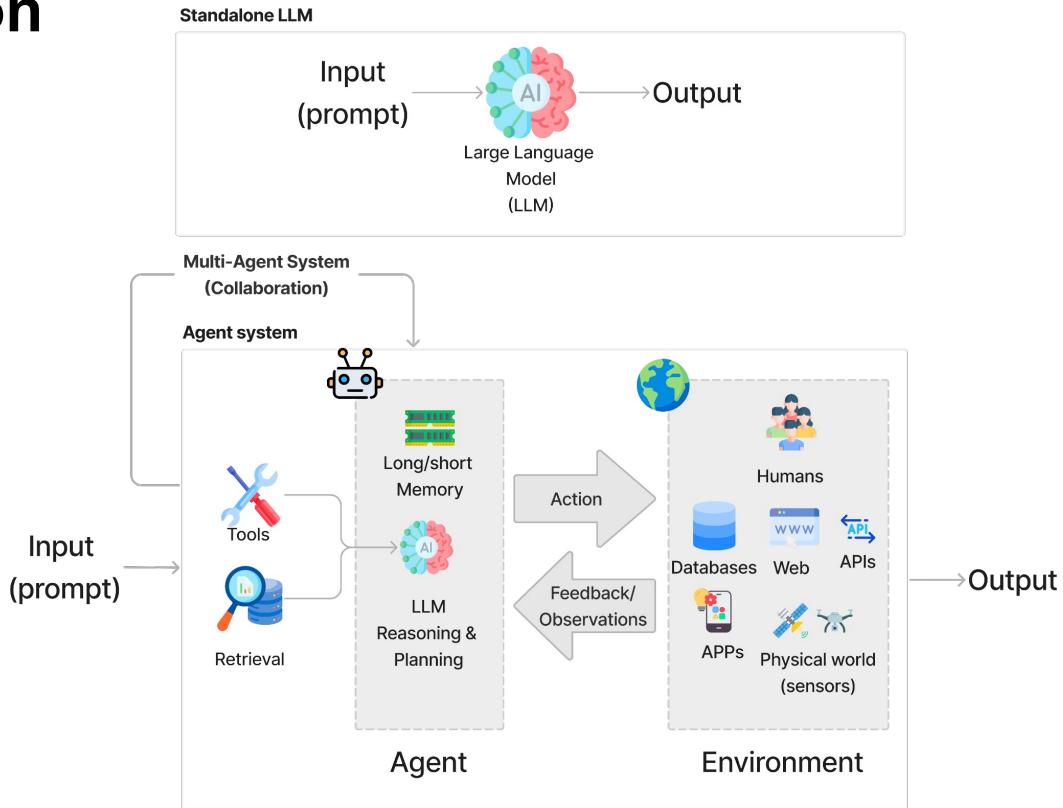


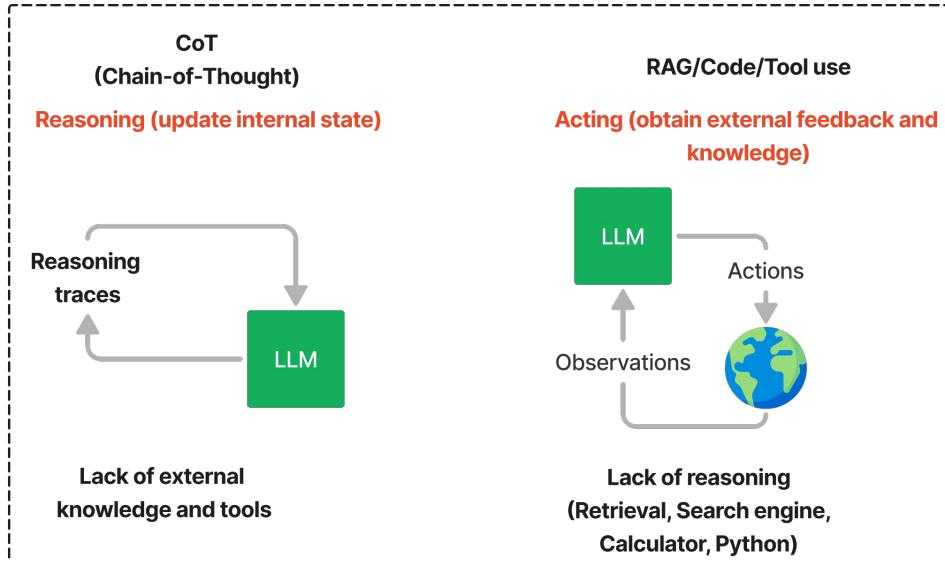
```
1  from dotenv import load_dotenv
2  from langchain import hub
3  from langchain.agents import (
4      AgentExecutor,
5      create_openai_tools_agent,
6  )
7  from langchain_openai import ChatOpenAI
8  from langchain_google_genai import ChatGoogleGenerativeAI
9  from tools import ImageCaptionTool, ObjectDetectionTool
10
11 load_dotenv()
12
13 tools = [ImageCaptionTool(), ObjectDetectionTool()]
14 # Get the instructions prompts for the agent
15 prompt = hub.pull("hwchase17/openai-tools-agent")
16 prompt.pretty_print()
17 # llm = ChatGoogleGenerativeAI(model="gemini-pro")
18 # llm = ChatOllama(model="llama2")
19 llm = ChatOpenAI(
20     temperature=0,
21     model_name="gpt-3.5-turbo"
22 )
23 # Construct the OpenAI Tools agent
24 agent = create_openai_tools_agent(llm, tools, prompt)
25 # Create an agent executor by passing in the agent and tools
26 agent_executor = AgentExecutor(agent=agent, tools=tools, verbose=True)
27 response = agent_executor.invoke({
28     "input": "describe what is in the image images/test.jpg, and return the
29     objects found",
30 })
31 print(response)
```

Agents (An unified solution for reasoning and Act)

A verbal agent is a service that interacts with a **Large Language Model (LLM)** to perform goal-oriented tasks by leveraging available tools and contextual information. It relies on the LLM's natural language understanding capabilities to support reasoning, planning, and decision-making.

Agents can autonomously execute complex tasks and solve multifaceted problems.



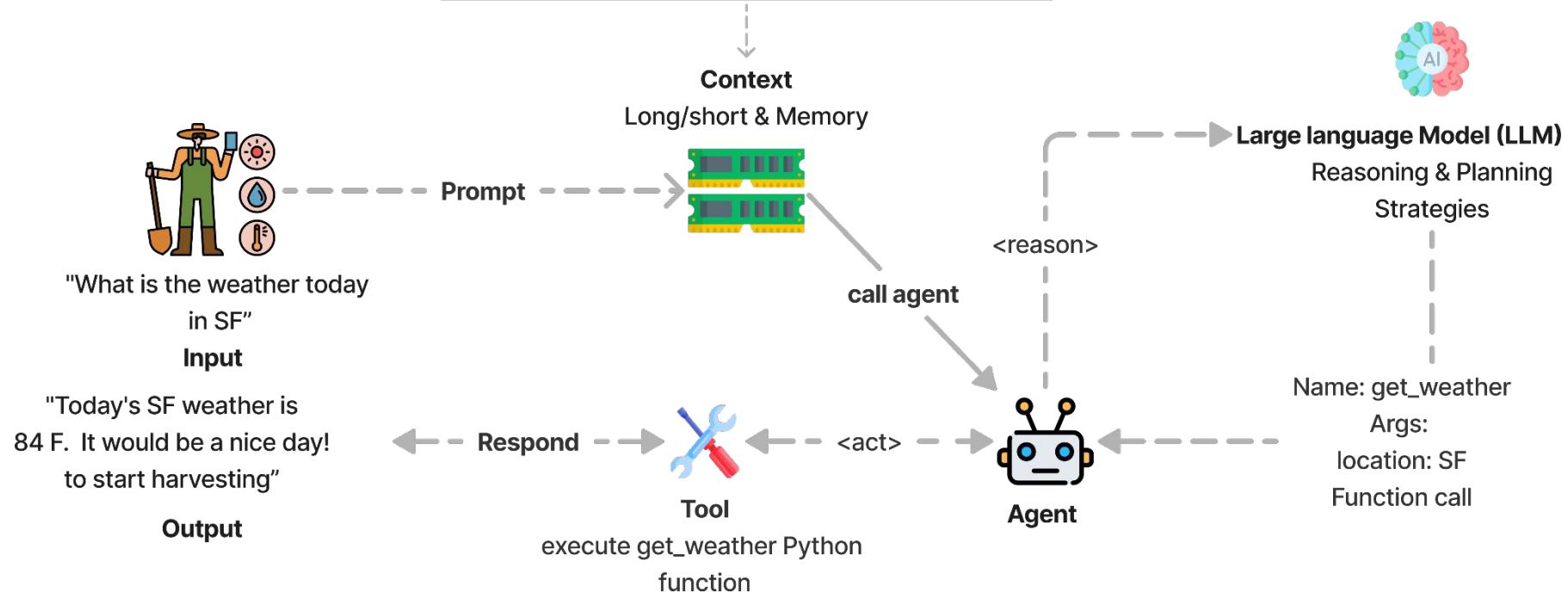


ReAct: A zero-shot emerging paradigm that introduces the concept of agents capable of both reasoning and taking actions

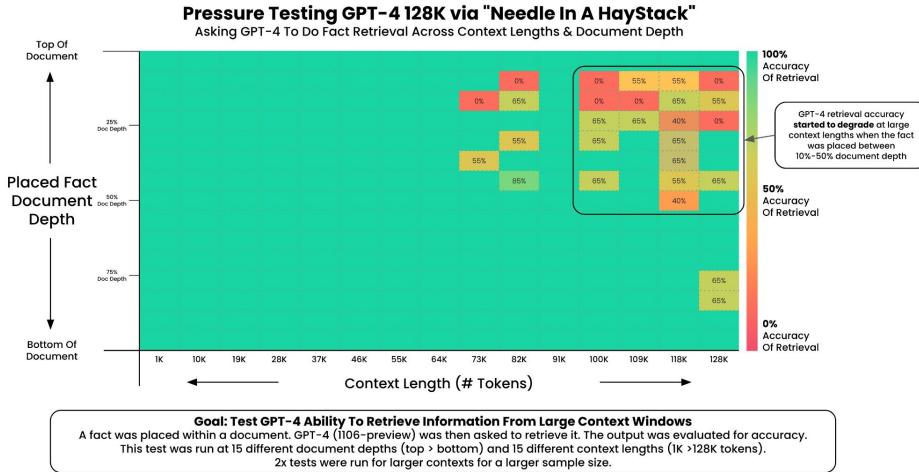


Agent Scratchpad

You are a helpful assistant with access to the following functions. Use them if required - {tools}. The output needs to be in the following format.....



Why do we need multi-agent architectures?



“Needle in a Haystack’ test

.... Interestingly, if the “needle” is positioned towards the beginning of the context, the model tends to overlook or “forget” it — whereas if it’s placed towards the end or as the very first sentence, the model’s performance remains solid....

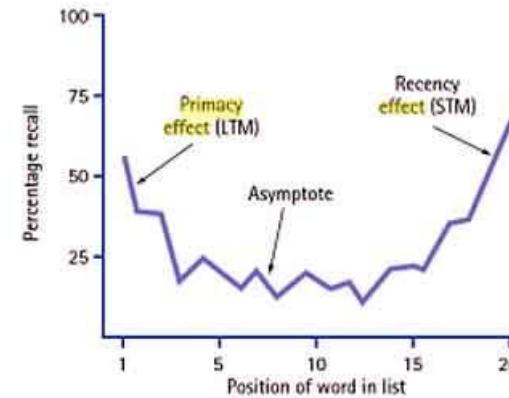
<https://arxiv.org/pdf/2307.03172>

<https://arxiv.org/pdf/1906.04341>

That phenomenon also occurred in humans.

The observation that LLMs tend to focus more on the beginning and end of a sequence has interesting parallels with biological systems, particularly in human cognition and memory. This effect can be likened to the **serial position effect** in psychology, which encompasses two phenomena:

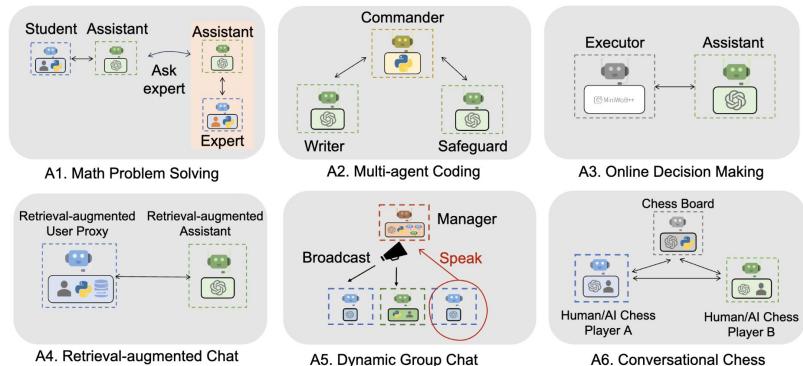
1. **Primacy effect** – People tend to remember the first items in a sequence more easily.
2. **Recency effect** – People tend to recall the most recent (or last) items more effectively.



Murdock discovered that recall probability depends on a word's position in a list, with better recall for words at the beginning (primacy effect) and end (recency effect), while middle words are often forgotten. Murdock's experiment involved a total of 103 participants, all of whom were students from an introductory psychology course.

Moving beyond single-agent systems to tackle complex tasks while mitigating and reducing hallucinations.

- A team of agents with diverse capabilities—language models and specialized tools—collaborating to solve complex tasks
-
- **Why?**
- Leverage a team of specialized agents (LLMs + tools)
- Group responsibilities to improve task success
- Focused agents outperform those managing many tools
- Use separate prompts for better performance
- Evaluate and improve agents independently without breaking the system



Agents collaboration patterns

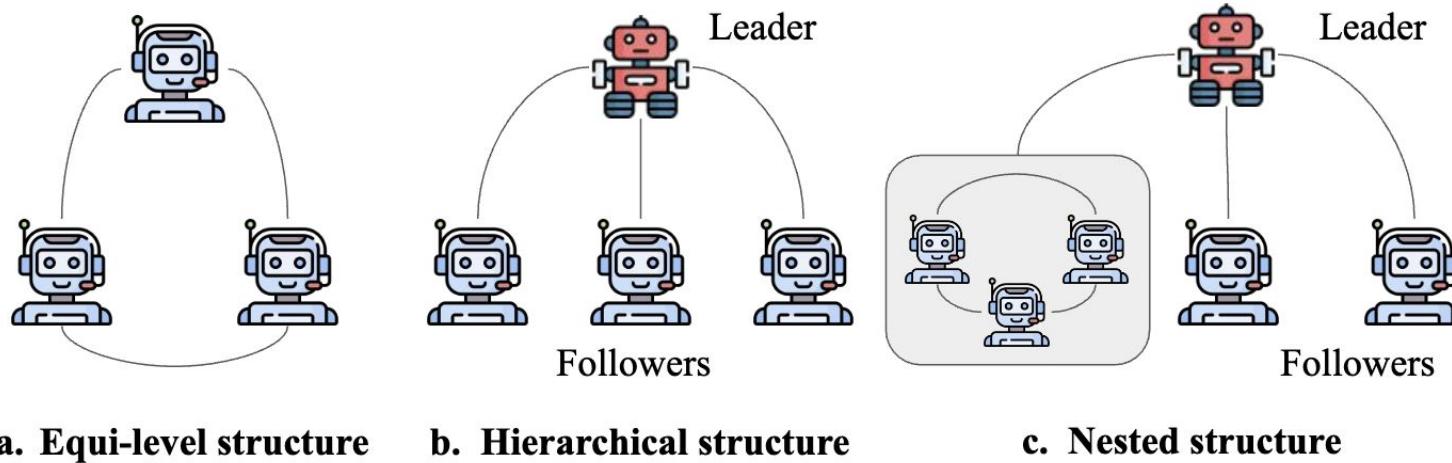


Figure 1. Structures of multi-agent systems.

Building Effective Agents

*"We've worked with dozens of teams building LLM agents across industries. Consistently, the most successful implementations use simple, composable patterns **rather than complex frameworks**."*

*"Agent" can be defined in several ways. Some customers define agents as 1) **fully autonomous systems** that operate independently over extended periods, **using various tools to accomplish complex tasks**. Others use the term to describe more 2) **prescriptive implementations that follow predefined workflows**. At Anthropic, we categorize all these variations as agentic systems, but draw an important architectural distinction between workflows and agents:*

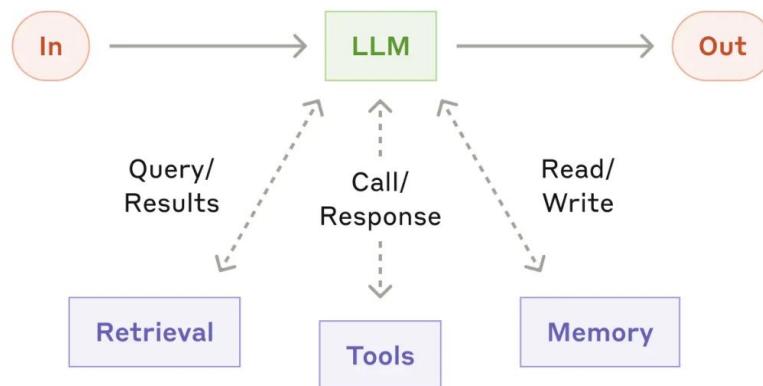
- **Workflows** are systems where LLMs and tools are orchestrated through predefined code paths.
- **Agents**, on the other hand, are systems where LLMs dynamically direct their own processes and tool usage, maintaining control over how they accomplish tasks.

– Anthropic

<https://www.anthropic.com/engineering/building-effective-agents>

Building block: The augmented LLM

The basic building block of agentic systems is an LLM enhanced with augmentations such as retrieval, tools, and memory.



Workflows

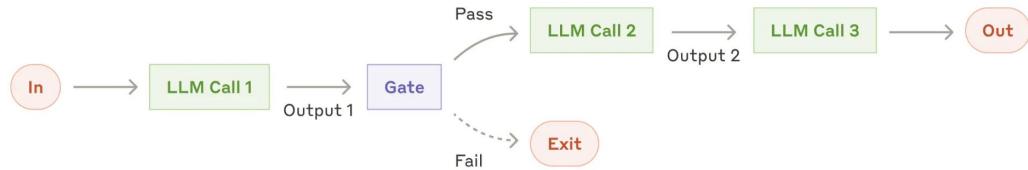
Workflow: Routing



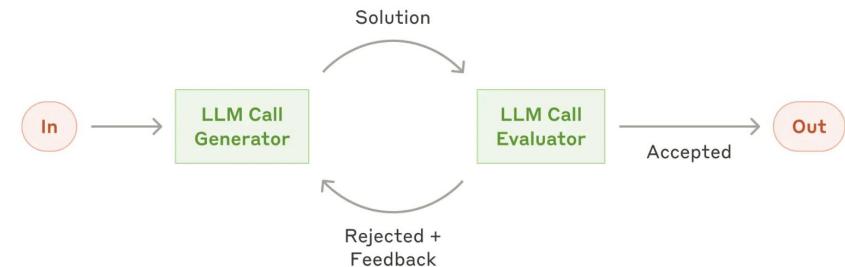
Workflow: Parallelization (Breaking a task into independent subtasks run in parallel & vote)



Workflow: Prompt chaining



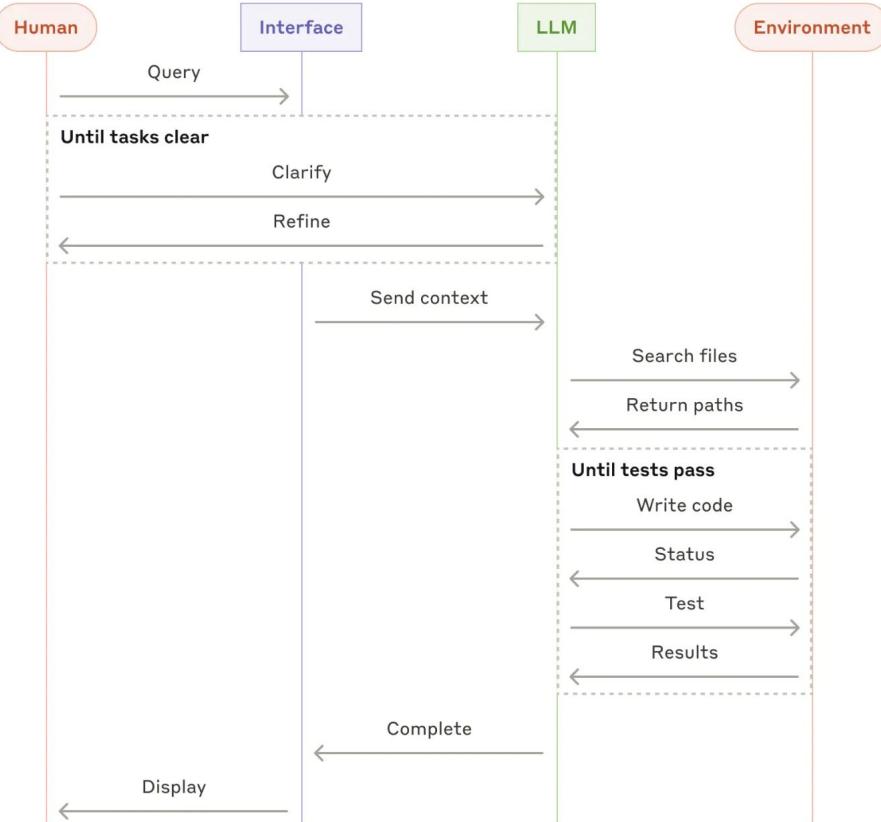
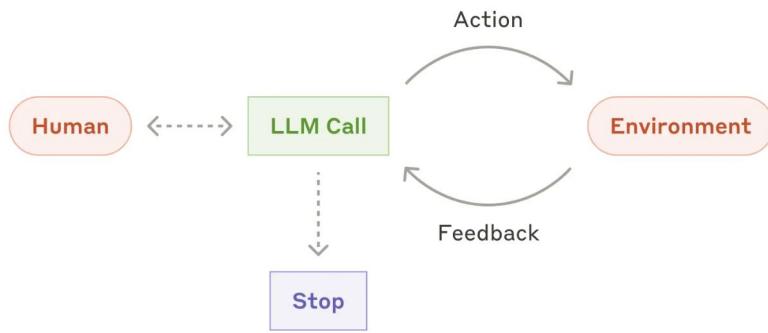
Workflow: Evaluator-optimizer (reflection)



Workflow: Orchestrator-workers (Supervisor)



Agents



Building a Multi-Agent System - *from scratch*

- Goal definition
- Baseline
- Tools
- Eval testbed
- Agent

This comes last,
not first!



Frameworks provide different Levels of Abstraction

Agent framework
(Agent Development Kit)



Low-level orchestration framework
(LangGraph)



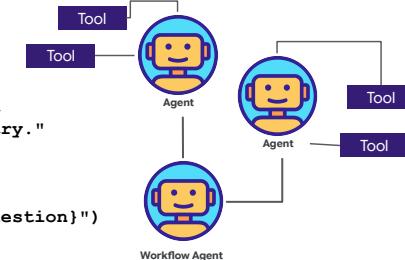
Low-level LLM framework
(Langchain)



DIY

```
capital_finder_agent = LlmAgent(  
    name="capital_finder_agent",  
    model="gemini-1.5-flash-latest",  
    instruction="You are a helpful assistant  
    that provides the capital city for a given country."  
)
```

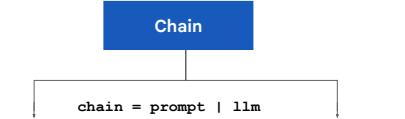
```
question = "What is the capital of France?"  
print(f"Calling ADK Agent with question: {question}")
```



```
workflow = StateGraph(GraphState)  
workflow.add_node("llm_call", call_llm)  
workflow.add_edge(START, "llm_call")  
workflow.add_edge("llm_call", END)
```



```
prompt = ChatPromptTemplate.from_messages([  
    {"role": "user", "content": "What is the capital of France?"}  
)  
chain = prompt | llm  
capital = chain.invoke({ "question": question })
```



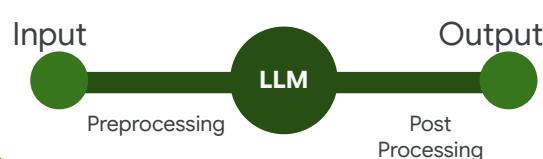
```
client = OpenAI()  
  
response = client.chat.completions.create(  
    model="gpt-3.5-turbo",  
    messages=[  
        {"role": "user",  
         "content": "What is the capital of France?"}  
    ]  
)
```

Vendor Specific

Restful

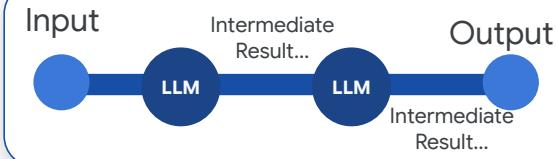
Architectural Patterns

Single Model Architecture



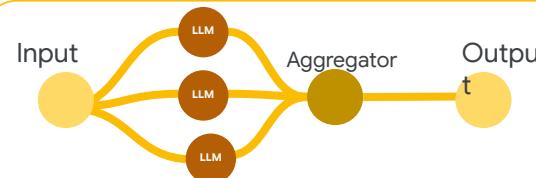
One LLM for all tasks

Sequential Models



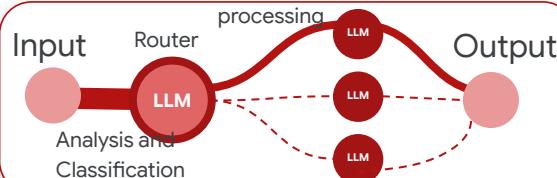
Chains multiple specialized LLMs in a pipeline.

Parallel Models



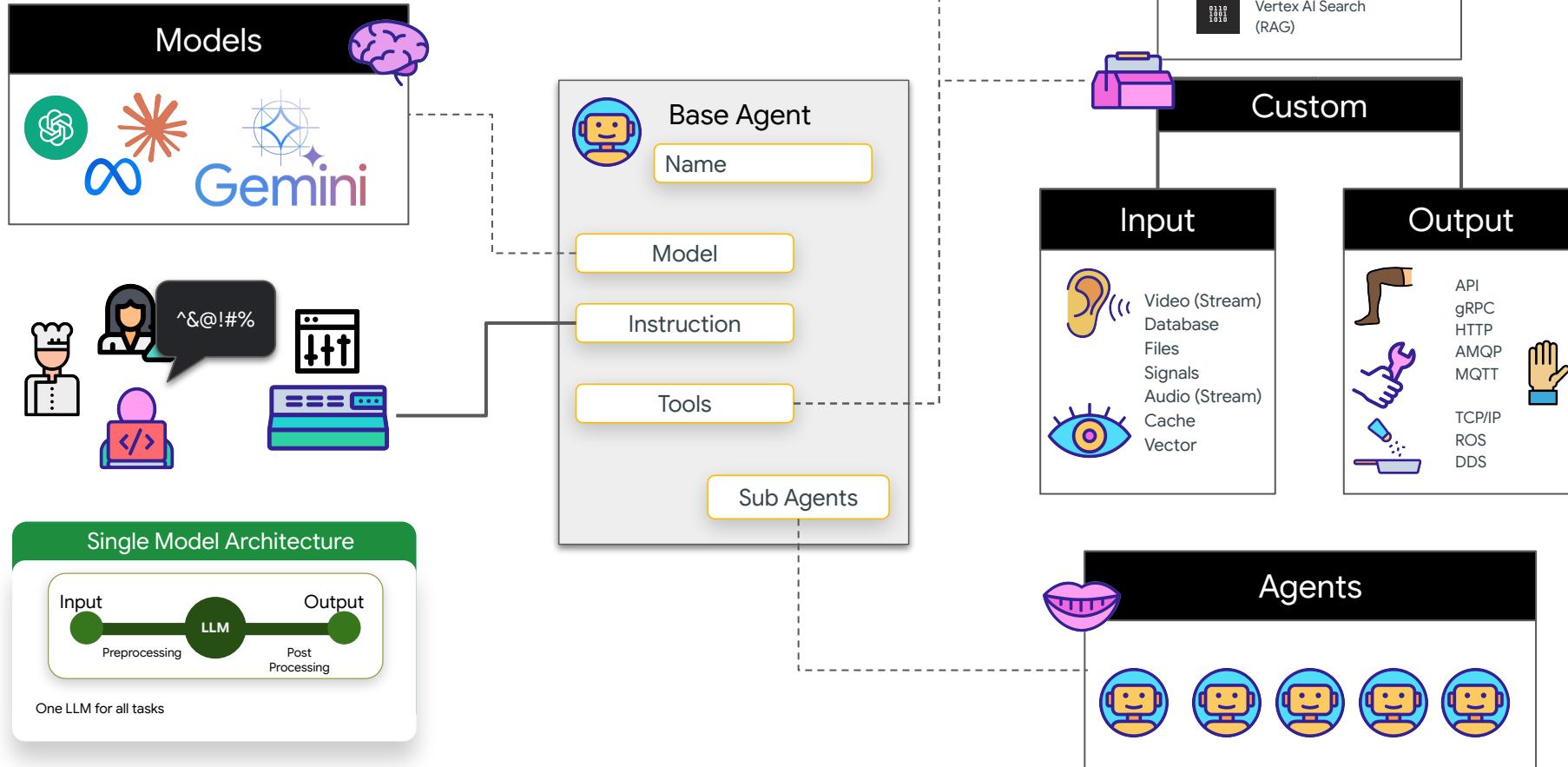
Multiple LLMs concurrently and aggregate outputs

Hierarchical Routing

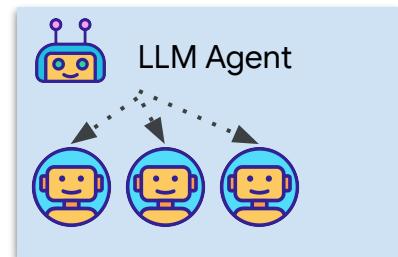
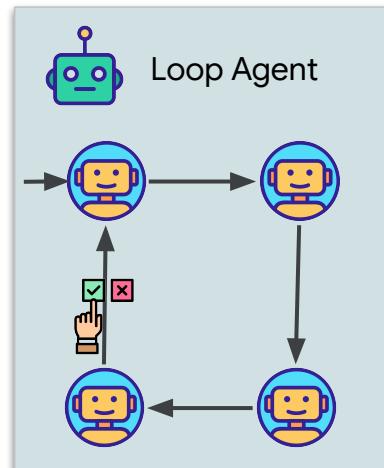
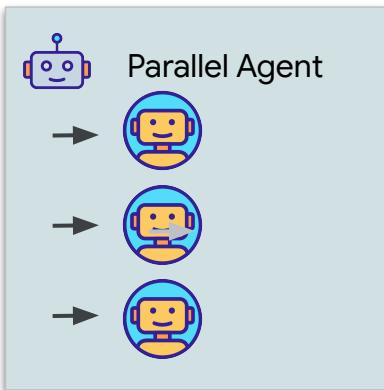
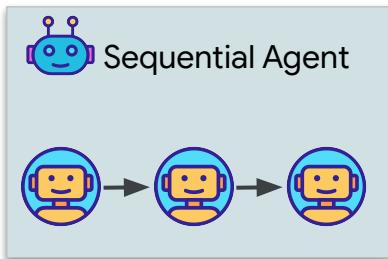
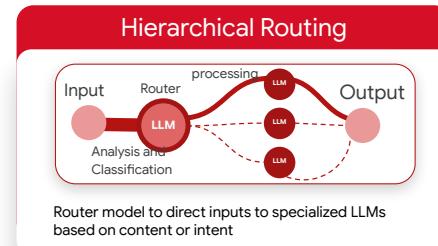
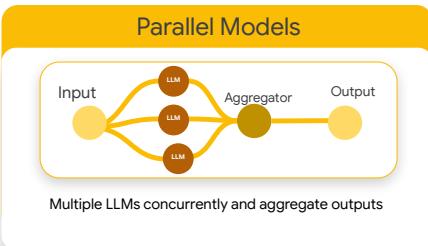
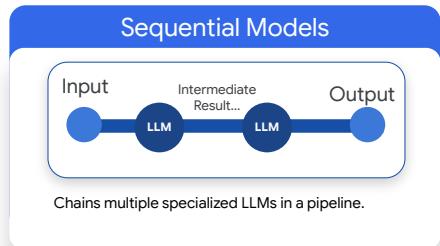


Router model to direct inputs to specialized LLMs based on content or intent

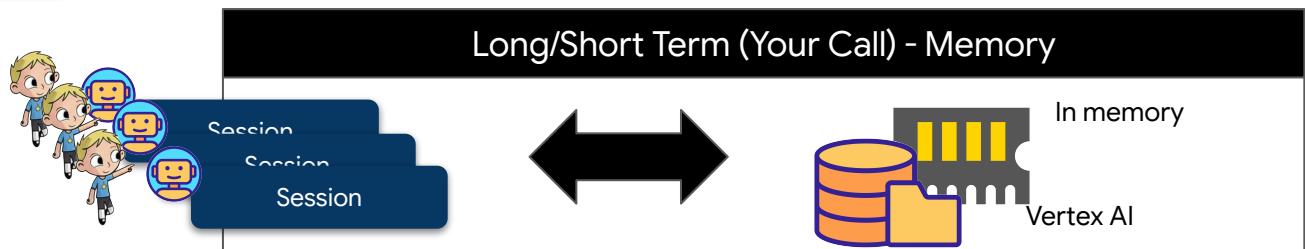
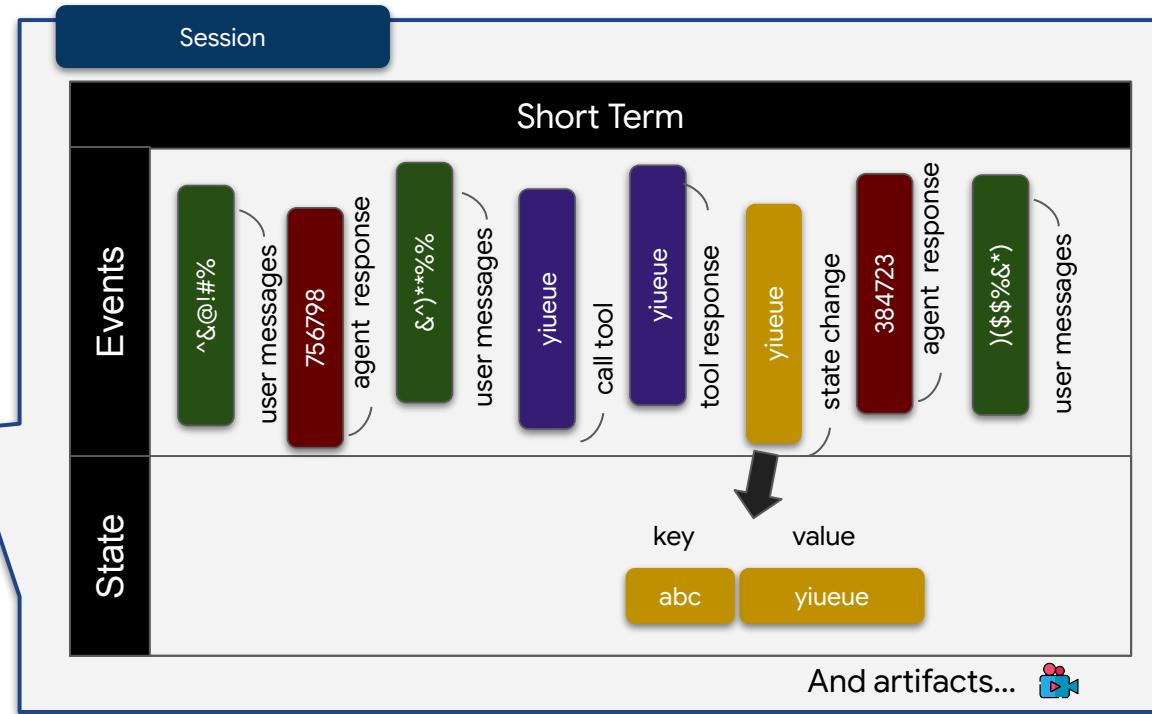
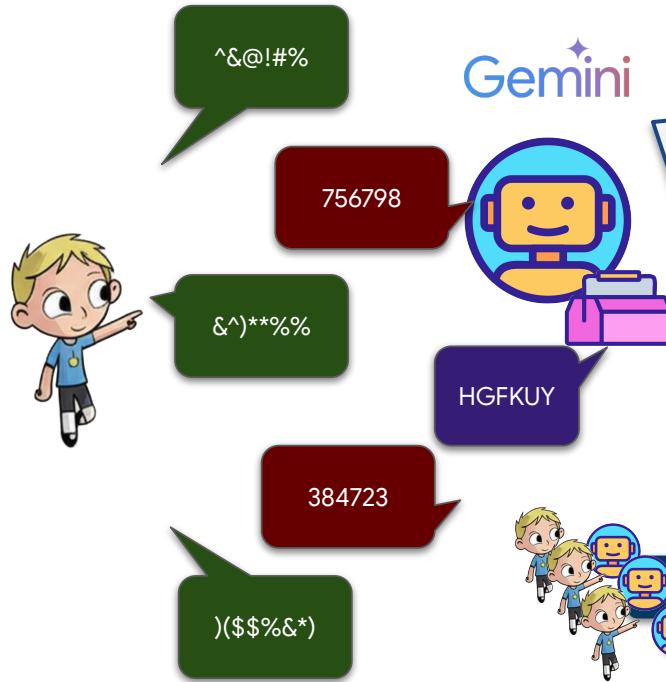
ADK Agent - Base



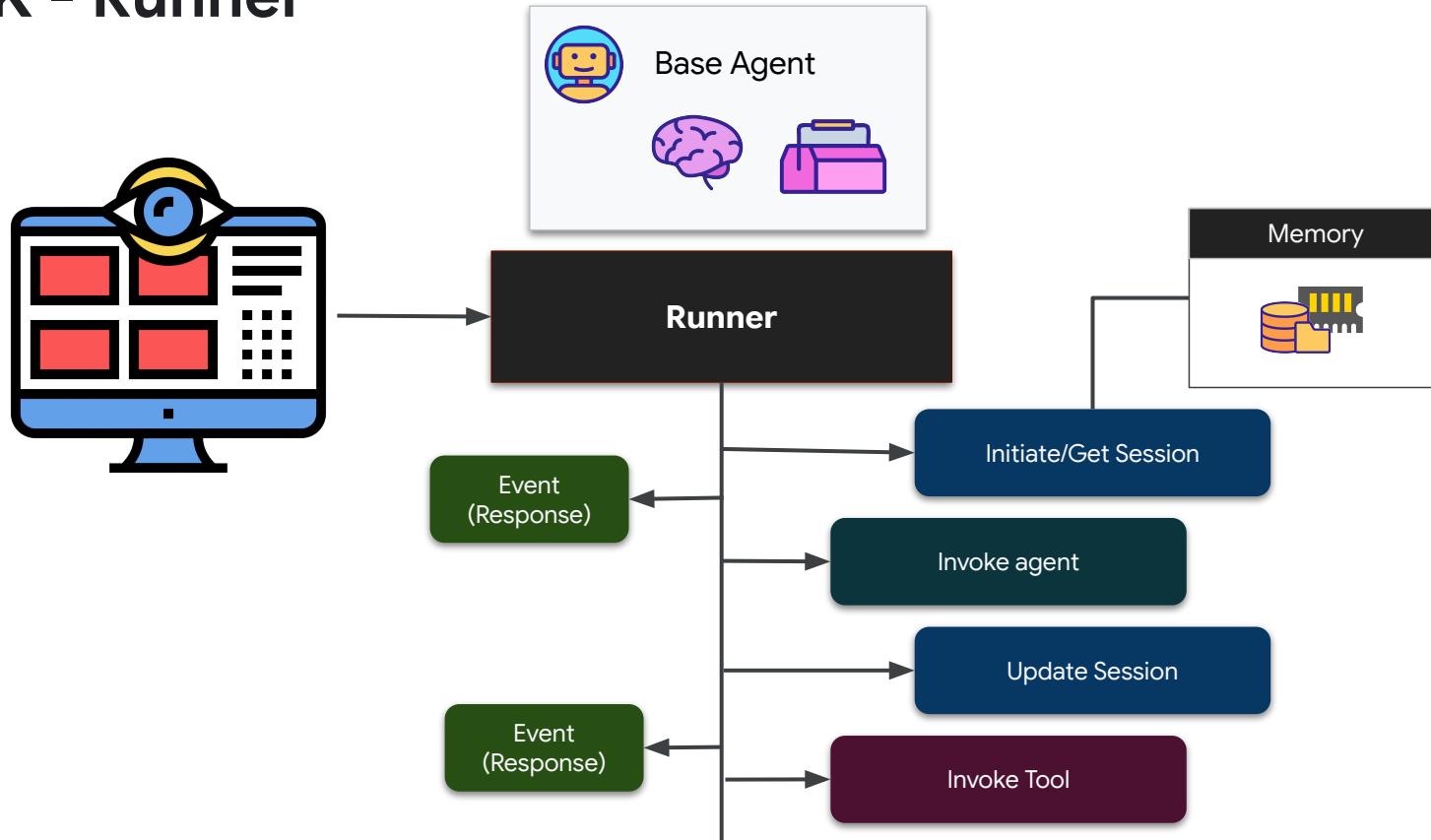
ADK Agent - Workflow



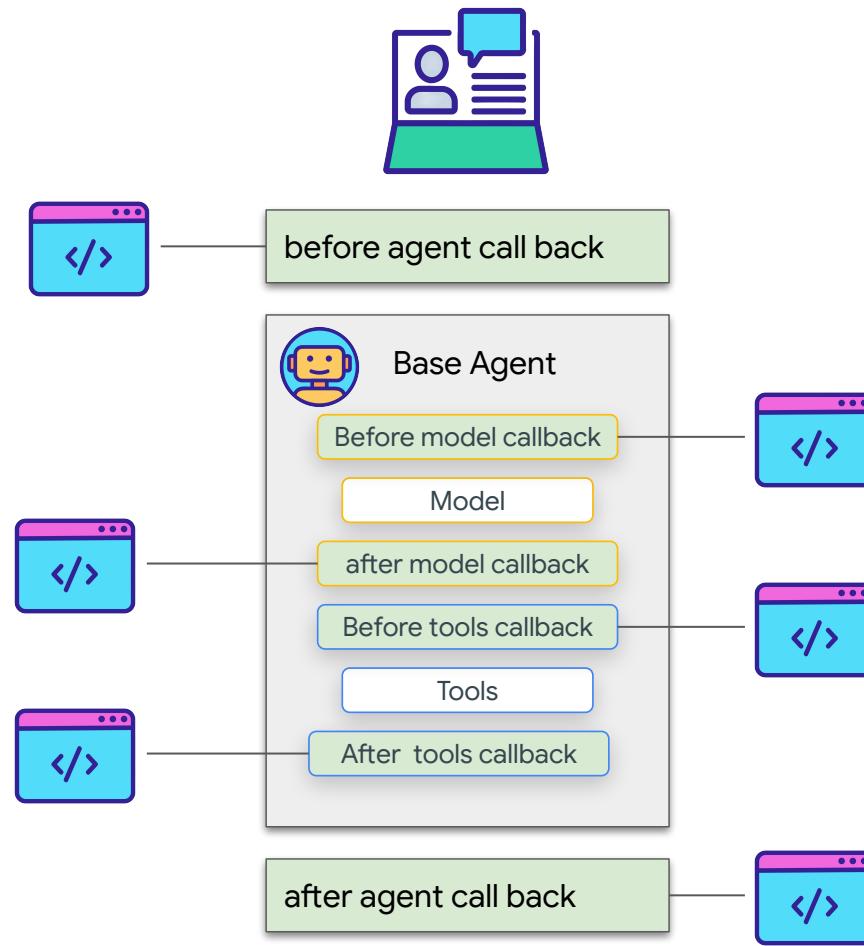
ADK - Memories



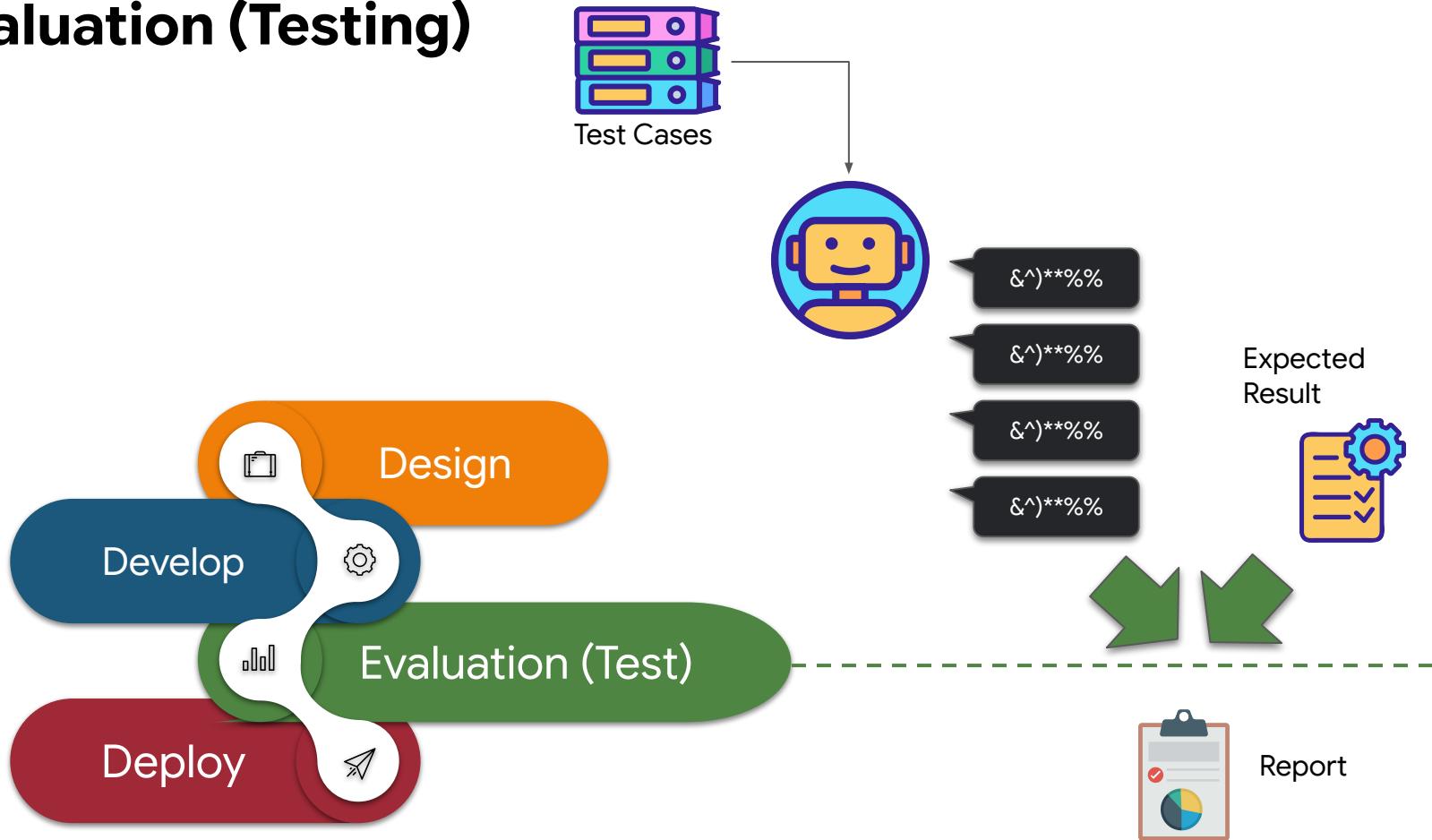
ADK - Runner



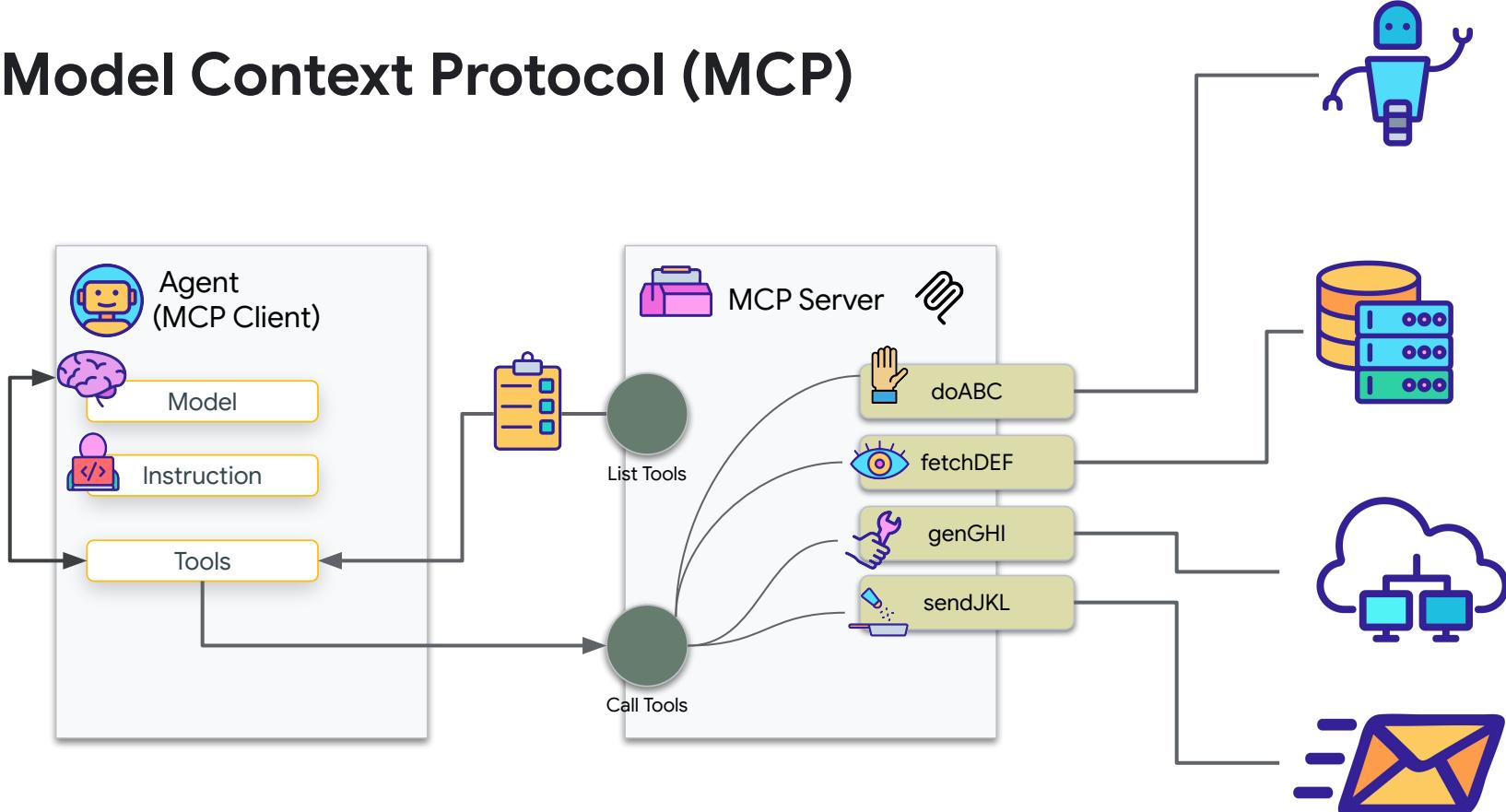
Call backs



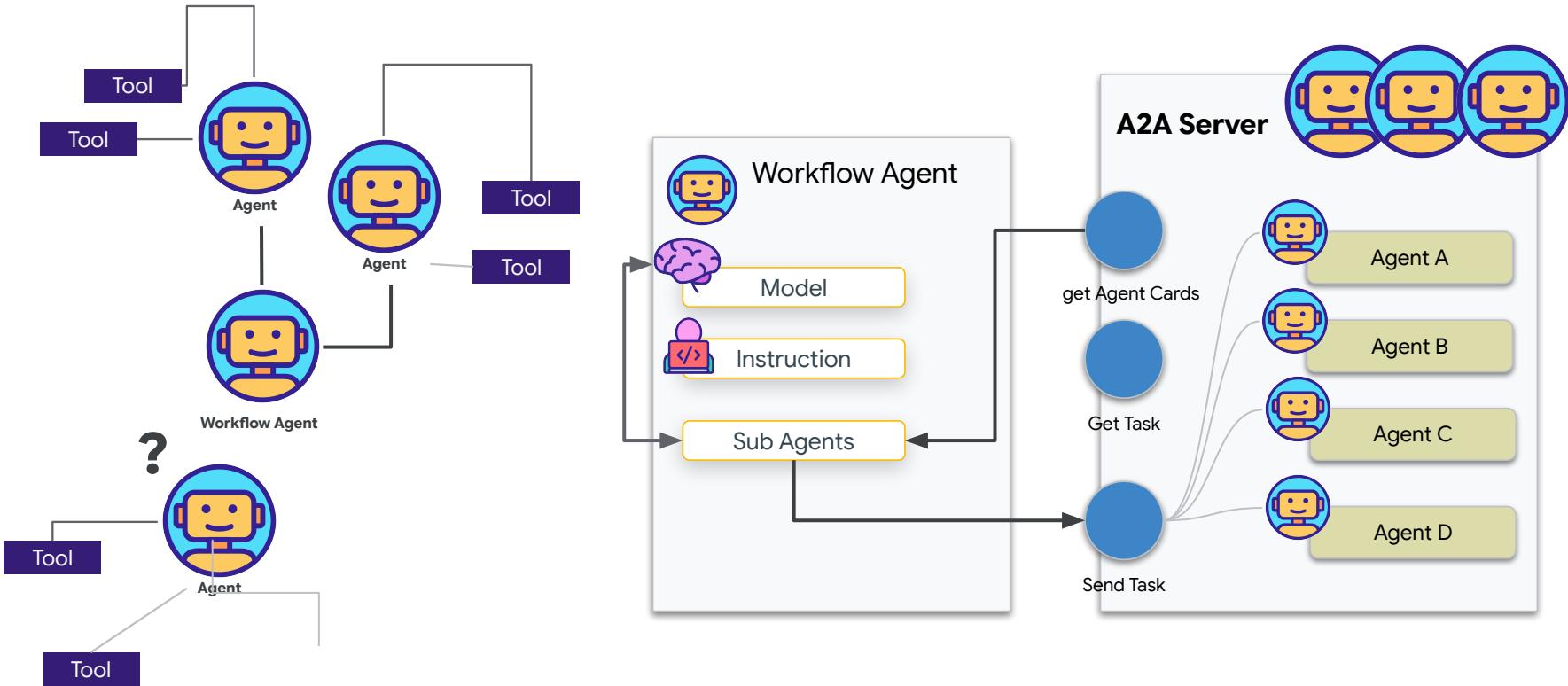
Evaluation (Testing)



Model Context Protocol (MCP)



Agent2Agent (A2A) Protocol



A2A - Agent Card

- A digital profile or manifest for an agent. It's a standardized metadata file (JSON)
- To orchestrating systems
- For discovery



Name: Agent ABC

Provider: Google (URL)

Preferred input/output: text

Capabilities:

- I can do streaming
- I can do Push Notifications

Skills:

- Skill A
 - Description
 - Input
 - Output
- Skill B
 - Description
 - Input
 - Output

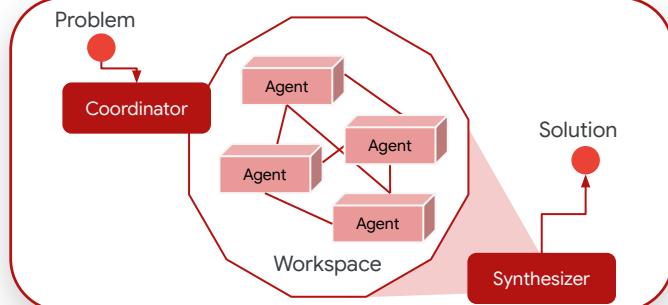
Authentication:

Scheme: Bearer

Key: xxx

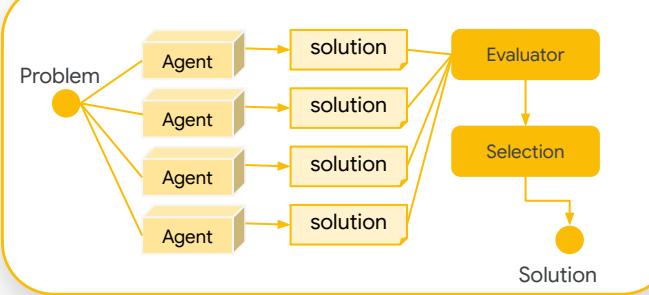
Multi-Agent Interaction

Cooperative Multi-Agent System



Multiple LLM-based agents collaborating to solve problems

Competitive Multi-Agent System

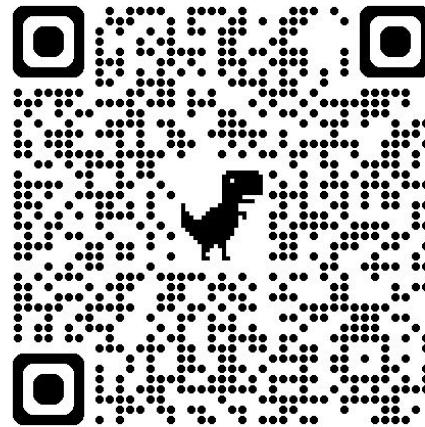


Independent LLM agents to propose competing solutions, with evaluator model selecting best or integrating elements

What is Next?

Advance Workshop: Building AI Agents

https://docs.google.com/document/d/1SVc_1-JVlo-6hsSUfF29jyiqT6vntRnzVz622TdjSyw/edit?tab=t.0





<https://rdi.berkeley.edu/events/agentic-ai-summit>

Google

Q&A