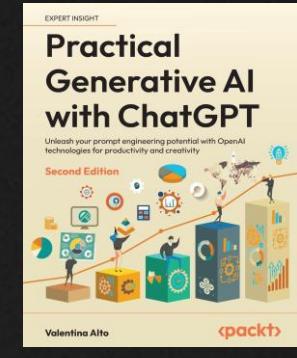




About Valentina Alto

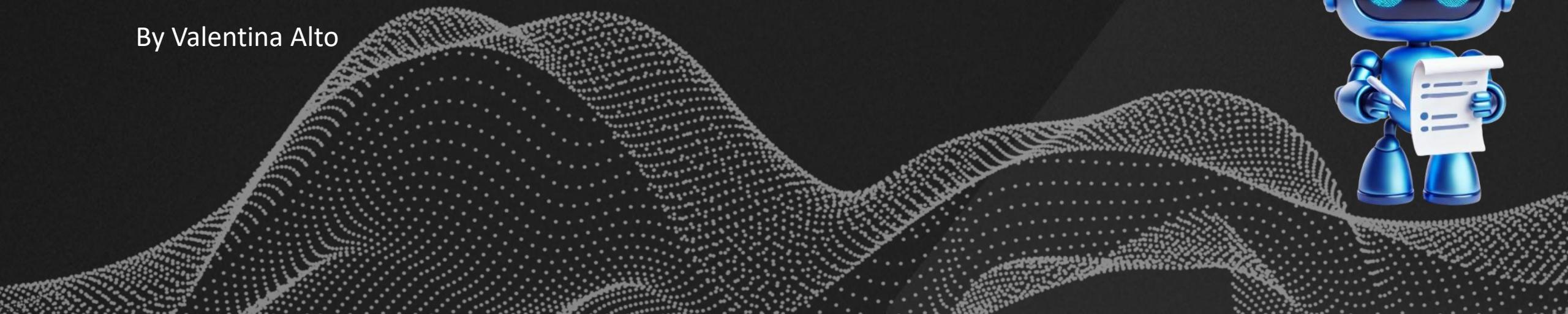
After completing her Bachelor's degree in Finance, Valentina Alto pursued a Master's degree in Data Science in 2021. She began her professional career at Microsoft as an Azure Solution Specialist, and since 2022, she has primarily focused on working with data and AI solutions in the Manufacturing and Pharmaceutical industries. Valentina collaborates closely with system integrators on customer projects, with a particular emphasis on deploying cloud architectures that incorporate modern data platforms, data mesh frameworks, and applications of Machine Learning and Artificial Intelligence. Alongside her academic journey, she has been actively writing technical articles on Statistics, Machine Learning, Deep Learning, and AI for various publications, driven by her passion for AI and Python programming.





AI AGENTS FOUNDATIONS – Day 1 | Part 2

By Valentina Alto



Agenda

-
- 01 Kick-Off and Introduction X

 - 02 From Content Generation to Agentic AI X

 - 03 Introduction to AI Agents X

 - 04 Demo Time X

 - 05 Introduction to Agentic Protocols: MCP and A2A X

 - 06 Conclusion X

01

Kick-Off and Introduction



Rapid Advancements of Generative AI



RAG

More context-aware, capable of maintaining coherent and relevant multi-turn conversations.



AI Agents

Generating sophisticated function or tool calls that can trigger actions in external systems



We are here!

Text generation

Creating human-like text, including articles, poetry, stories, and even computer code



Multi-modality

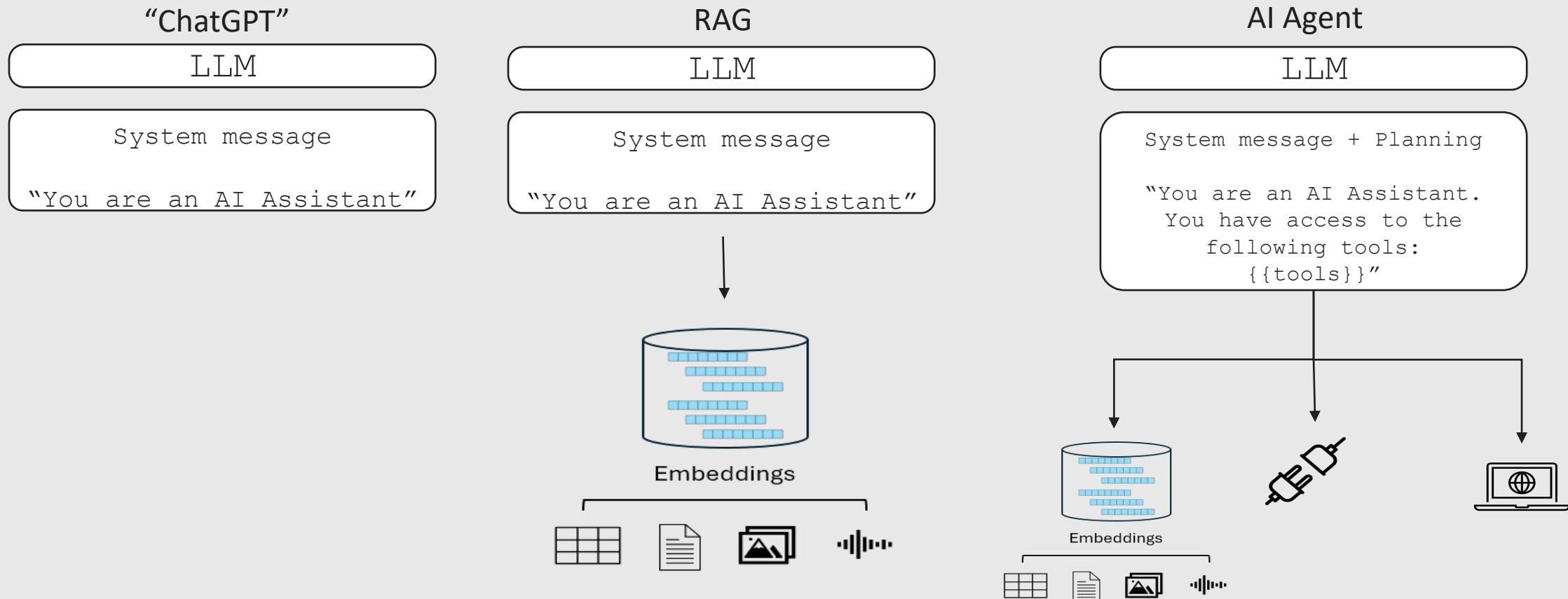
Understanding and generating data in different format including images, audio and video.



Multi-agents

Letting multiple AI agents interacting among each other in an autonomous way

AI Workflows

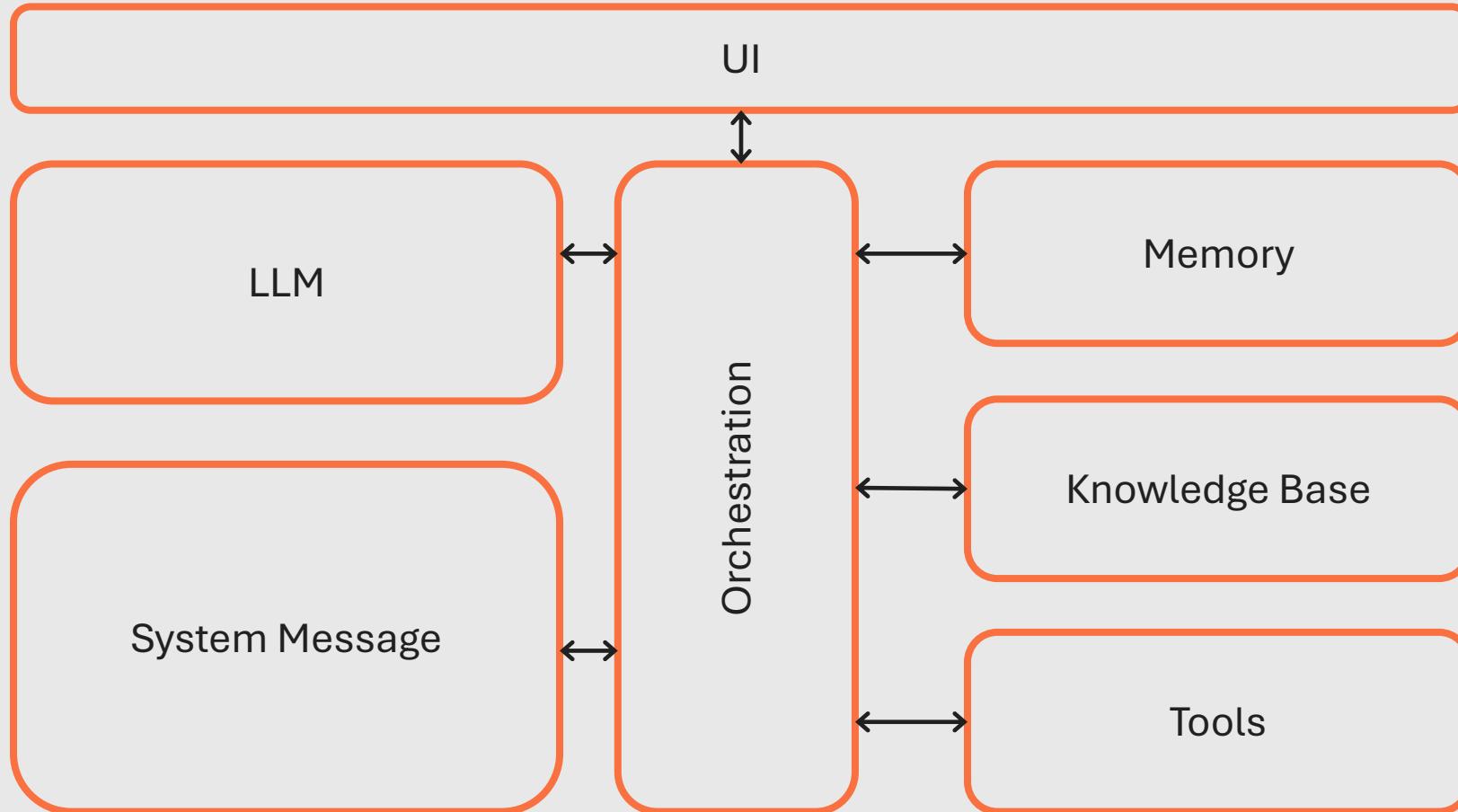


02

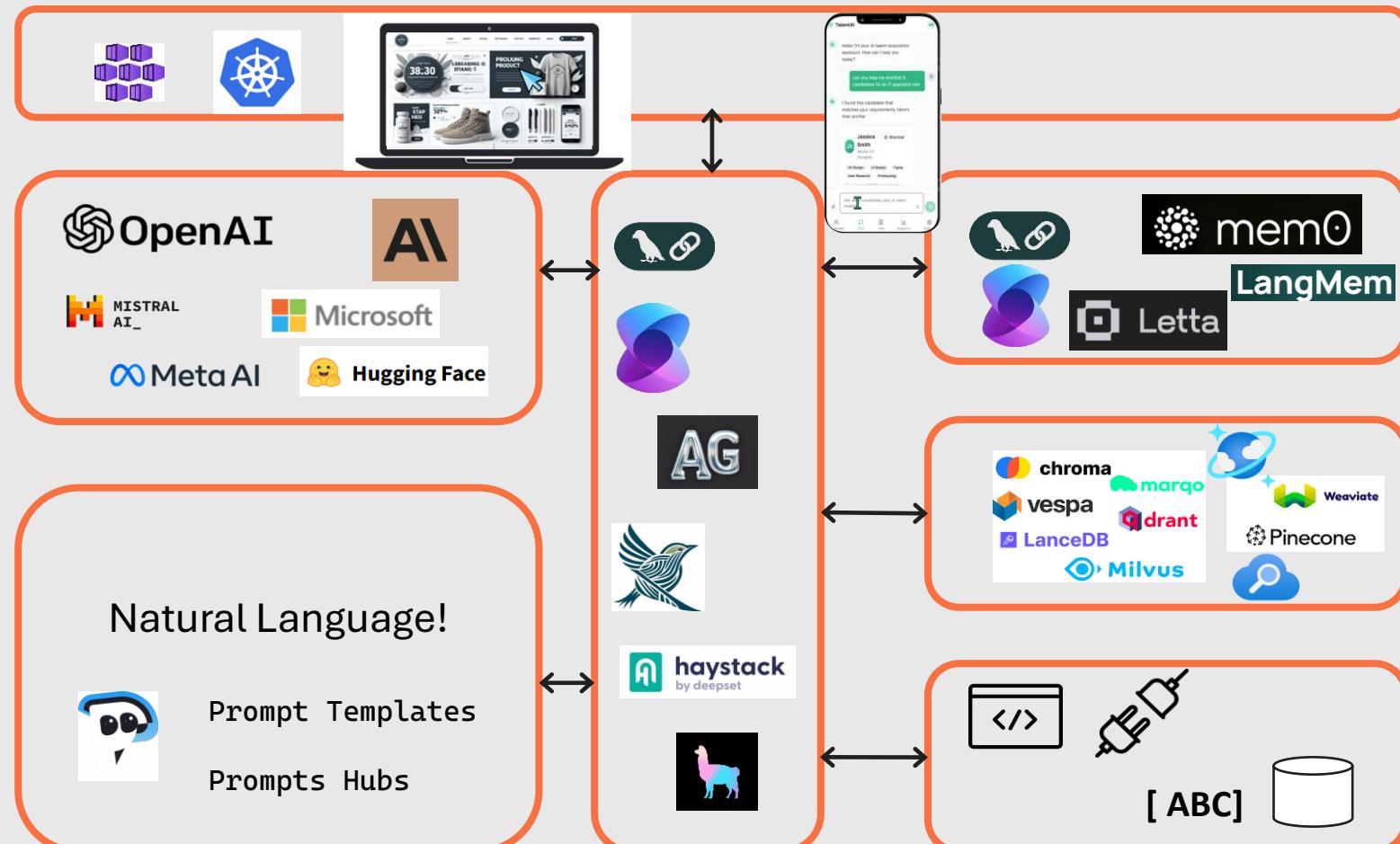
What is an agent?



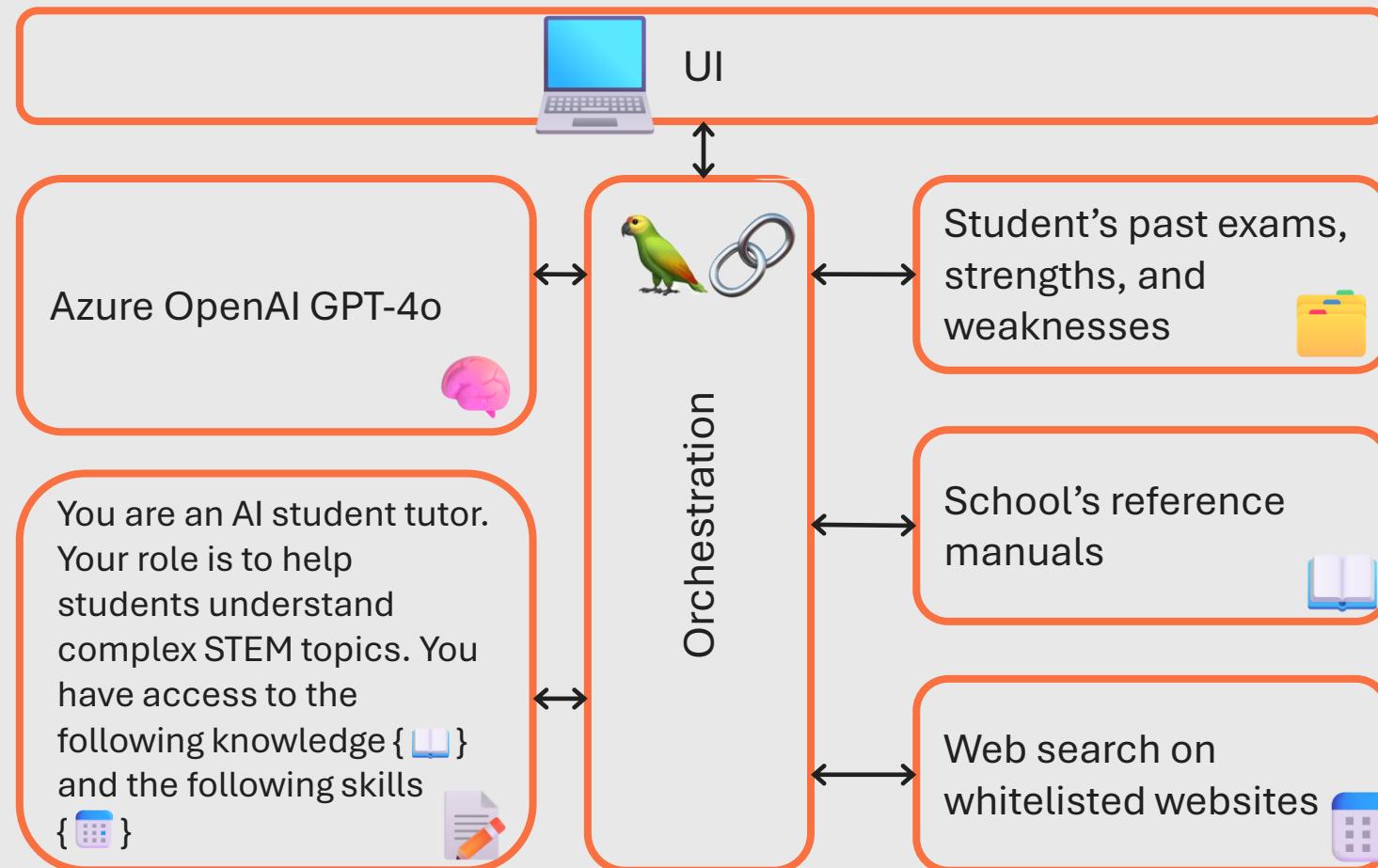
Anatomy of an AI Agent



Anatomy of an AI Agent



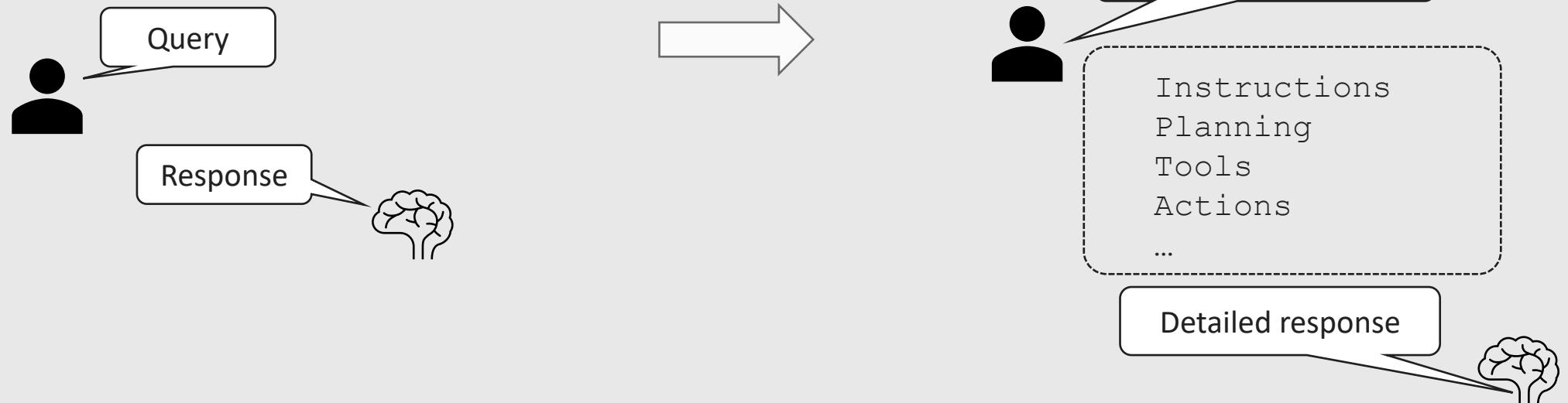
Example: AI-Powered Student Tutor



Paradigm Shift

From instantaneous, generic conversations...

...to highly specialized, complex tasks including autonomous decisions



Architecture



What's the weather tomorrow?

3



Metric

The forecast for tomorrow is sunny in Dubai, with an average temperature of 30 C.

Task Decomposition

Plan

I need to invoke the weather tool

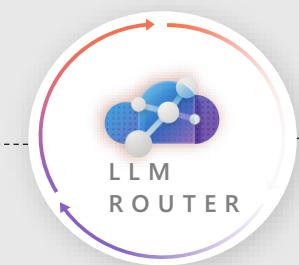
I need the “city” and “metric” parameters. To retrieve the city, I can sue the location tool.

To retrieve the unit, I’ll ask the user

I can now invoke the weather tool

Now I have the answer

You are a helpful AI Assistant



1

4

Weather Tool

`get_weather (city, unit)`

This is a tool to get the weather

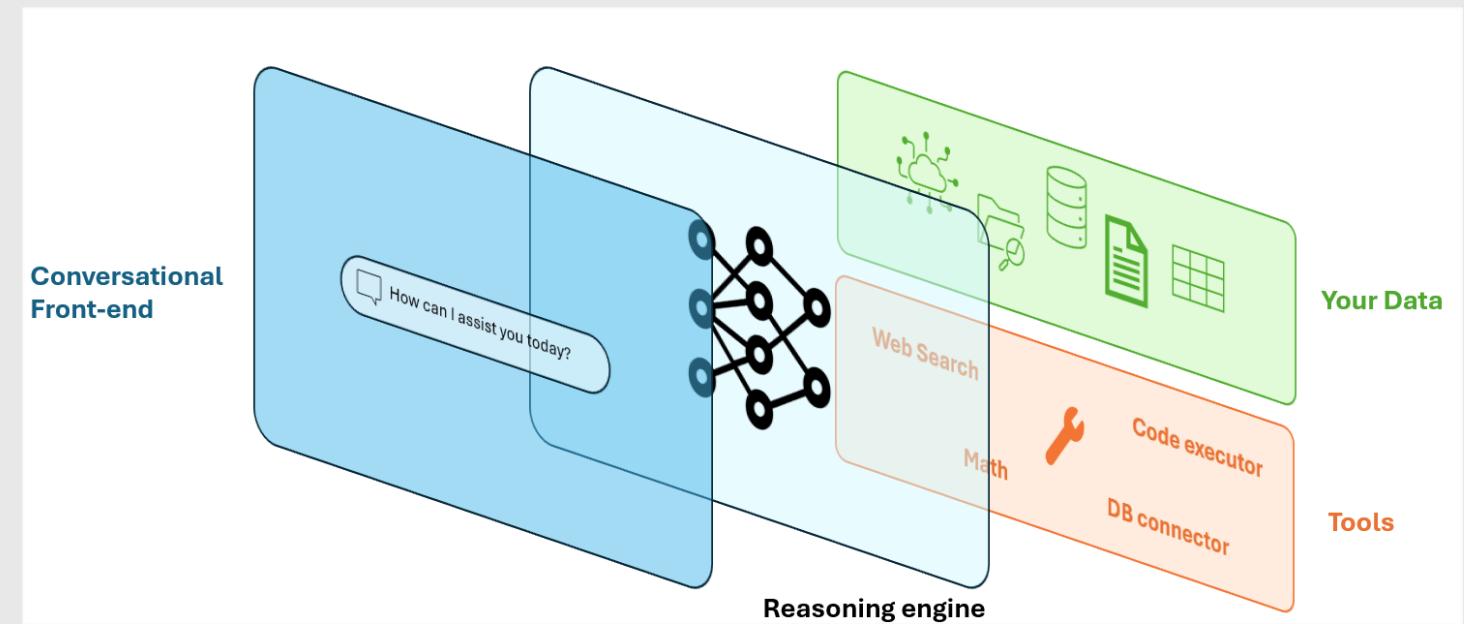
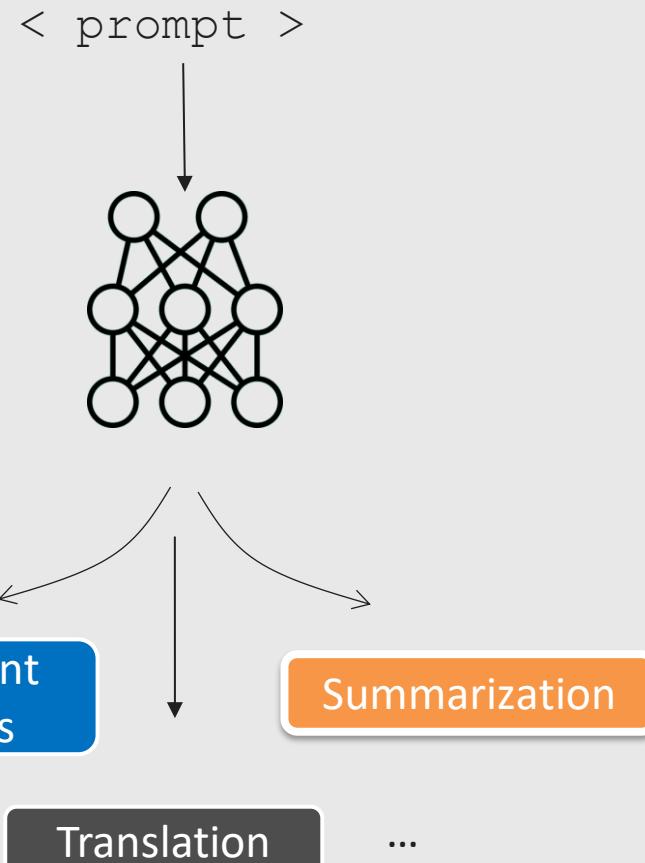
2

Location Tool

`get_location ()`

This is a tool to get the current location

Large Language Model



System Message

You are a helpful Civil Law Legal assistant. You receive queries from taxpayers. Refer to the Italian civil law system.

Always answer, providing a reference to the Codice Civile.
If the answer refers to another Legal domain, do not answer.

Always refer to the user with polite and empathetic style, making sure not to use an over-complicated language.

Before giving the response, review it and assess whether it might be potentially misleading. In that case, respond "I cannot help with that".

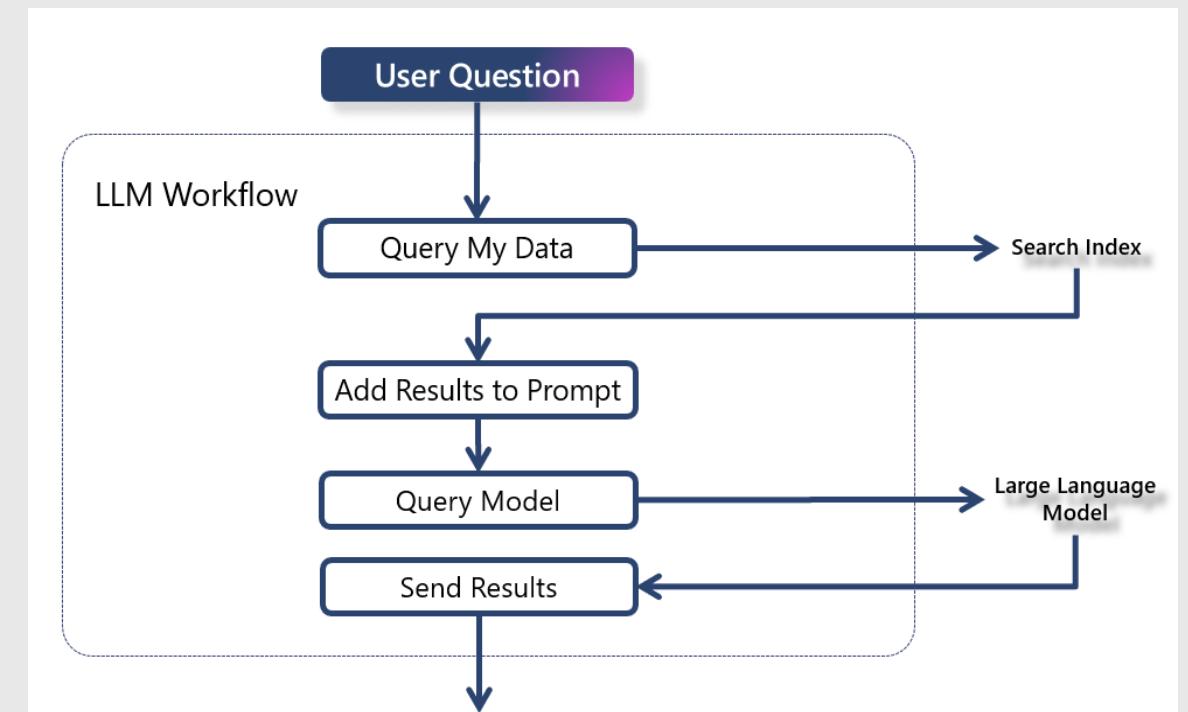
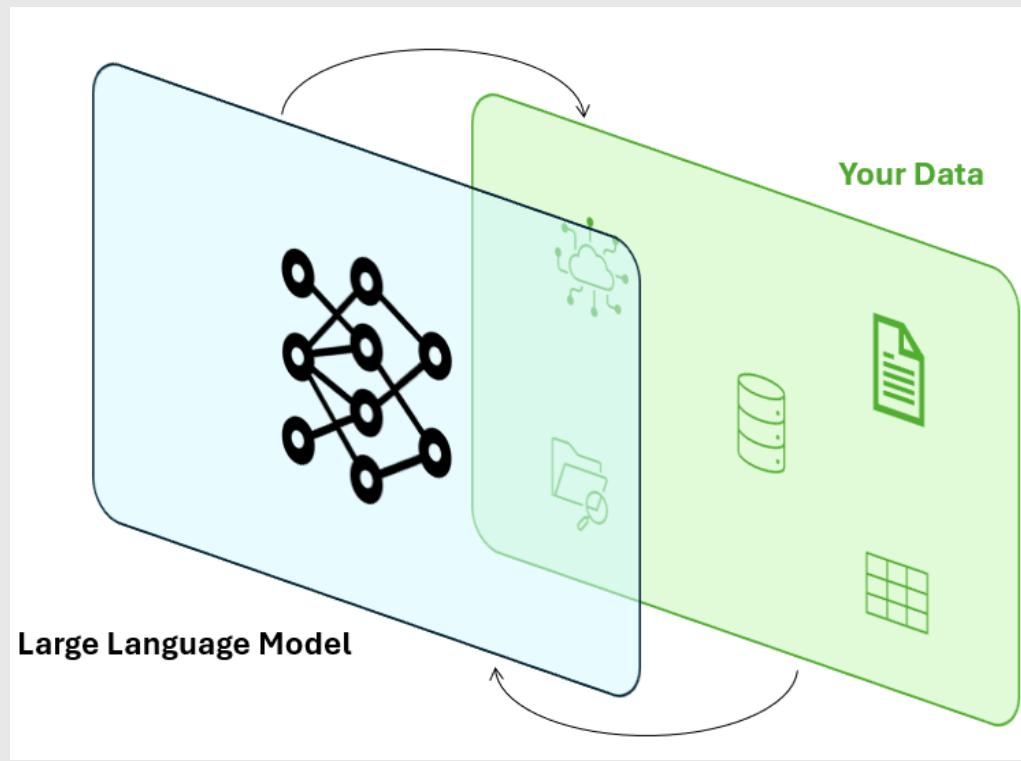
Setting the goal and scope

Grounding the model to relevant documents and preventing hallucination

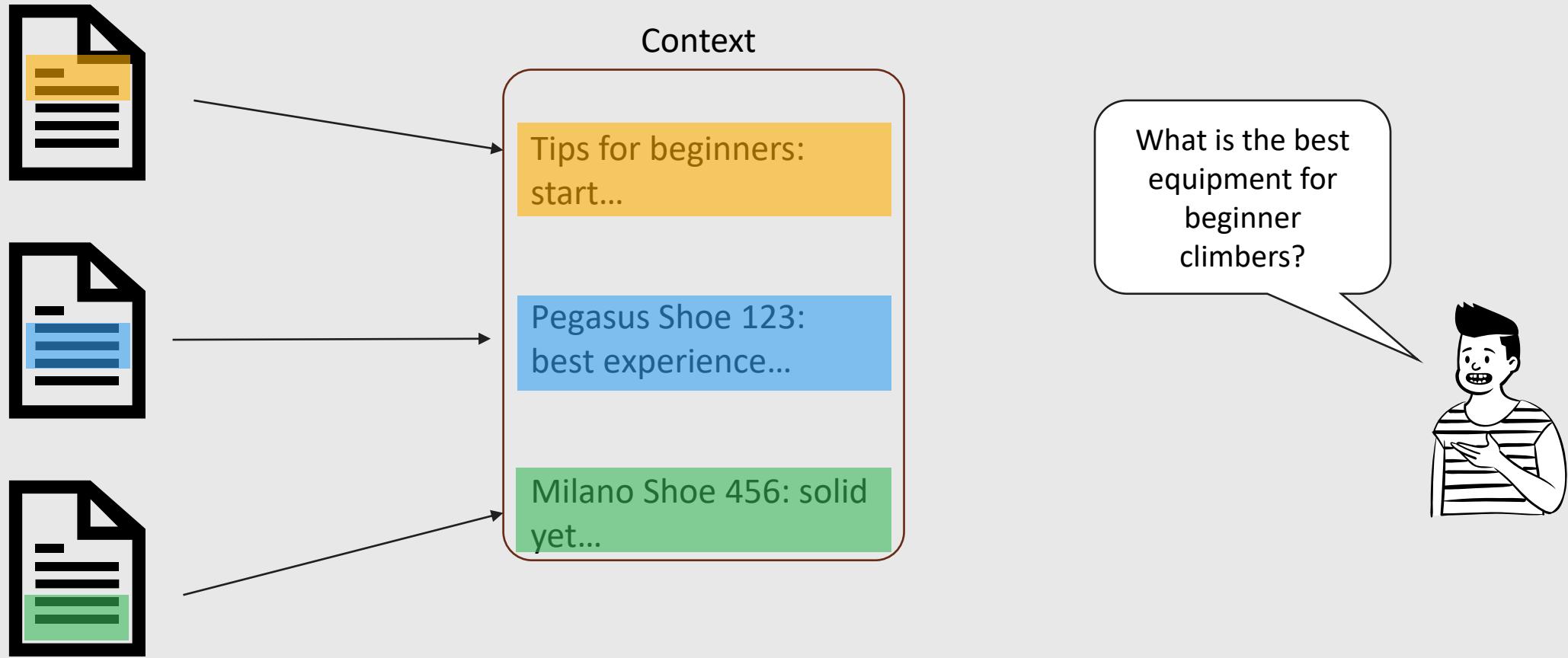
Setting the style of conversation

Incorporating responsible AI features

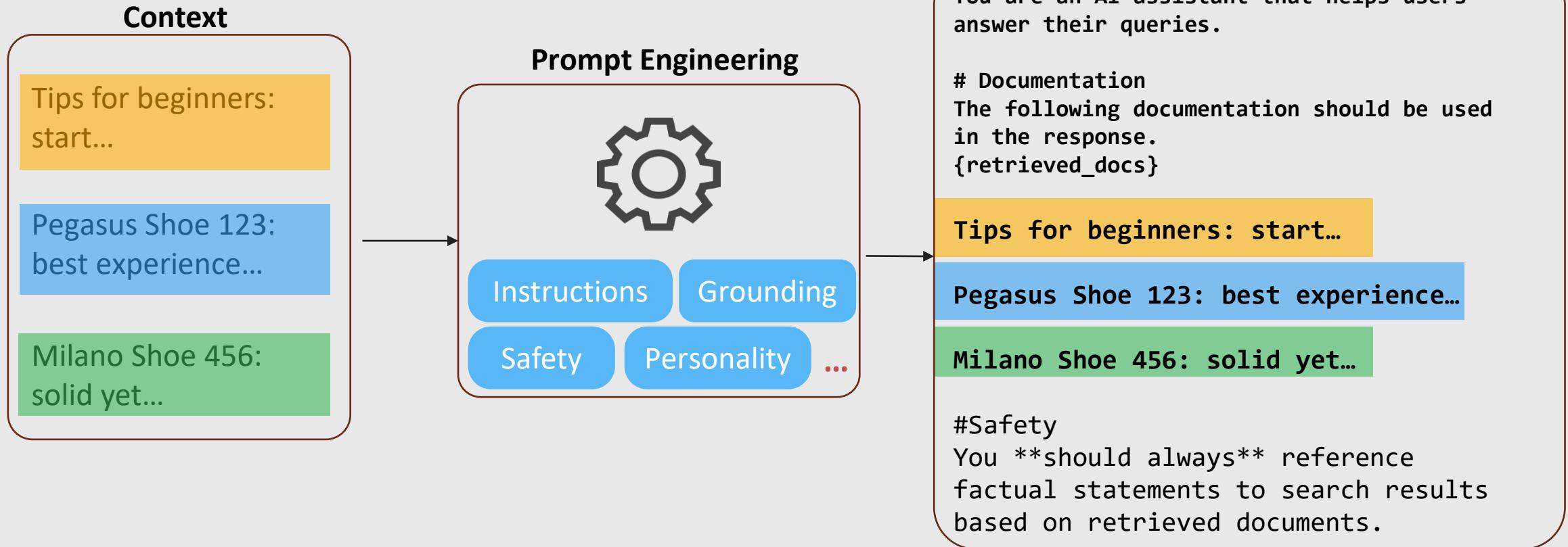
Knowledge - 1/8



Knowledge - 2/8



Knowledge – 3/8



Knowledge - 4/8

System Message

You are an AI assistant that helps users answer their queries.

Documentation

The following documentation should be used in the response.

{retrieved_docs}

Tips for beginners: start...

Pegasus Shoe 123: best experience...

Milano Shoe 456: solid yet...

#Safety

You ****should always**** reference factual statements to search results based on retrieved documents.

User's query

What is the best equipment for beginner climbers?

Generative Model
(e.g. GPT-4)

“According to the catalogue, if you are about to start climbing...”



Knowledge - 5/8

“cat”



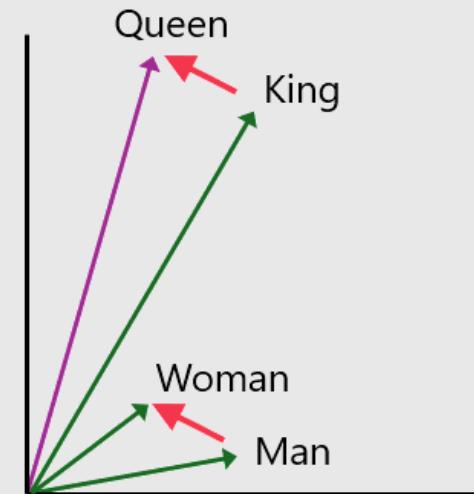
Embedding
model



{ 0.8 , 0.9 , -0.3 , -0.2 }

An embedding is a way of representing high-dimensional, non-numeric data, such as words or sentences, in a lower-dimensional space, such as vector.

A text embedding can capture the semantic and syntactic features of the text, such as meaning, context, and similarity.



King-Man+Woman ≈ Queen

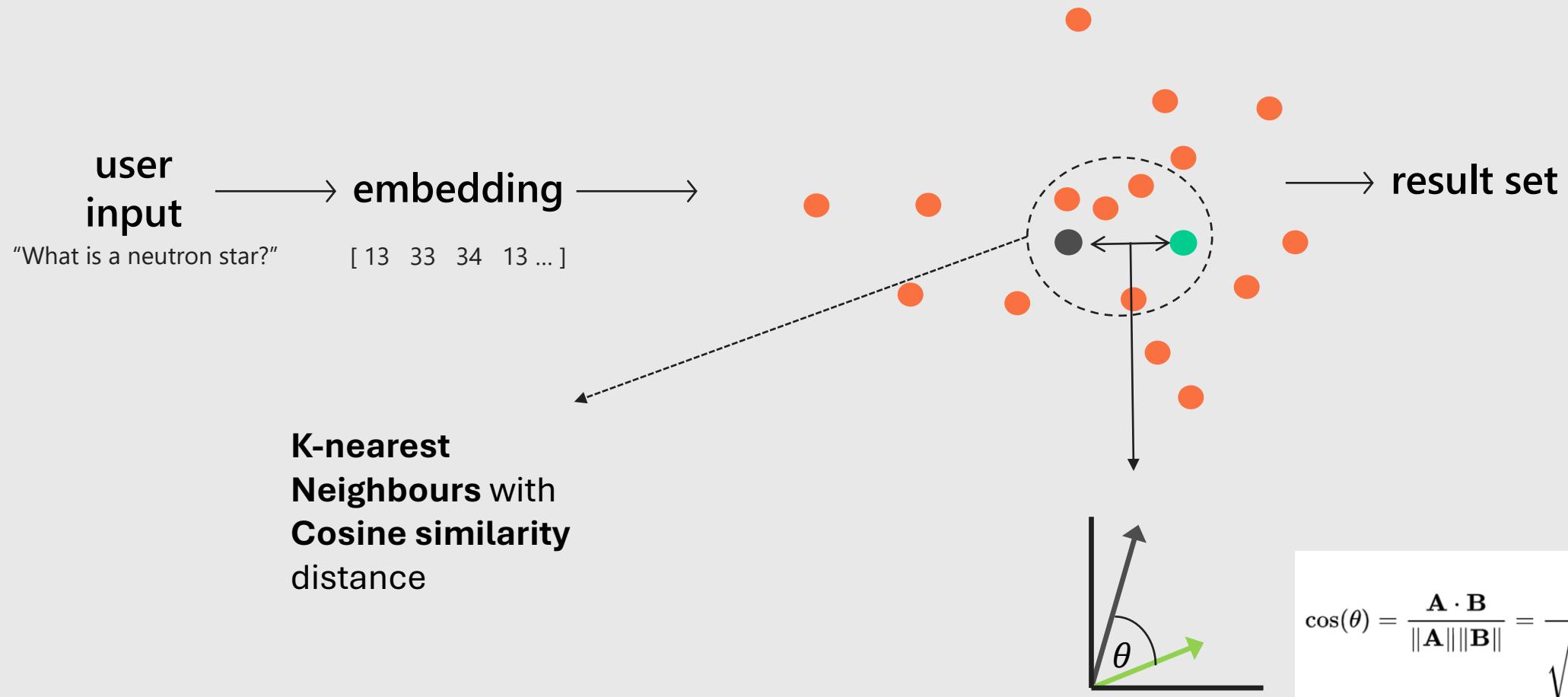
Knowledge – 6/8

Agent	Transcription
A	“Hello I’m calling for...”
B	“Am I speaking with ...”

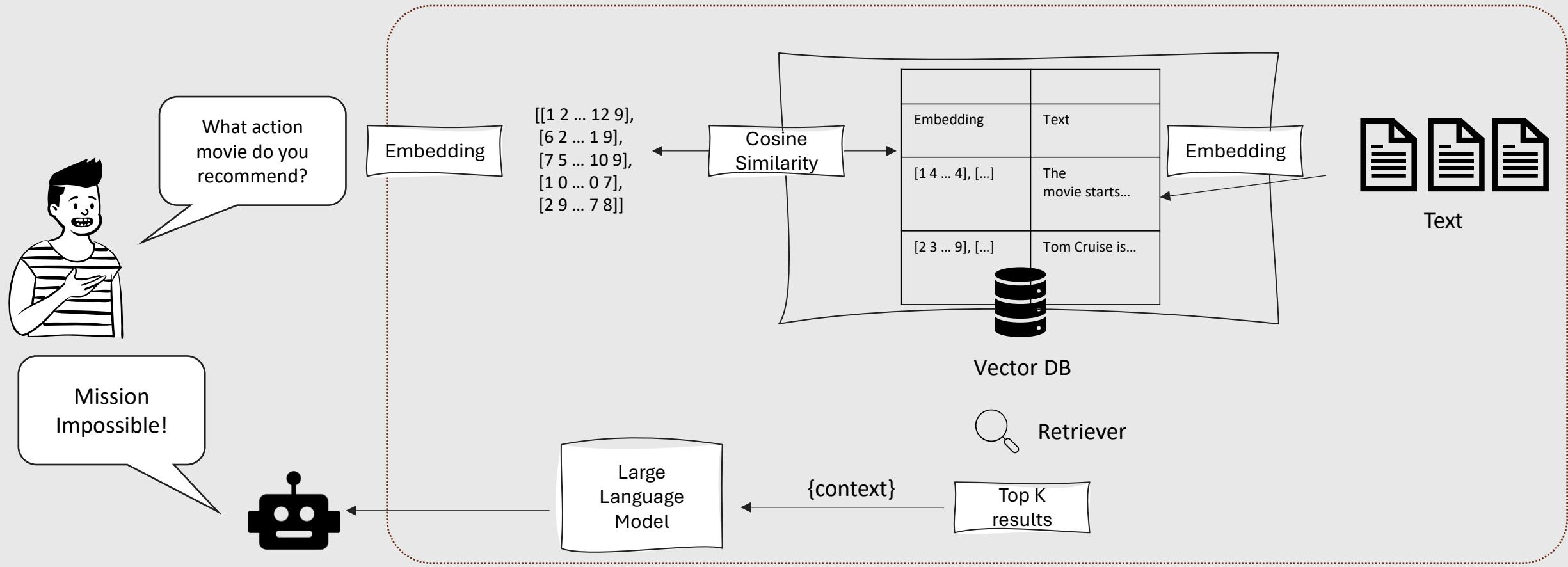
→ Embedding →

Agent	Transcription	Embedding
A	“Hello I’m calling for...”	[1 4 ... 4], [...]
B	“Am I speaking with...”	[2 3 ... 9], [...]

Knowledge - 7/8



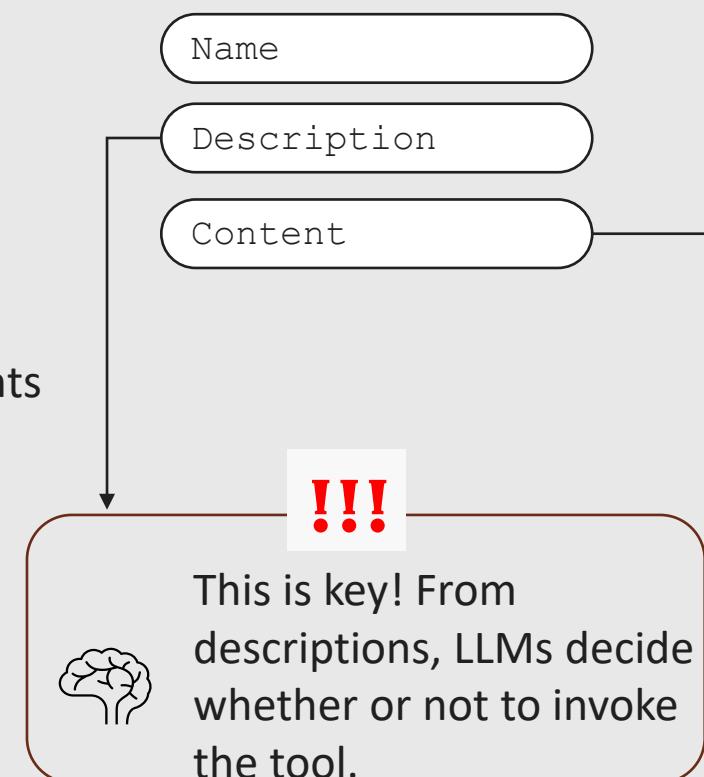
Knowledge - 8/8



Tools – 1/3



Tools allows AI Agents to interact with and take actions in the external world.



Pre-built connectors

```
import CRMConnector  
Connector = CRMConnector(client_id, api_key, ...)
```



Custom functions

```
def get_weather(location, metric):  
    ...  
    return ...
```



Semantic Skills

```
class SemanticSkills:  
    @kernel_function(name= "summarizer",  
                     description="summarize a list of docs")
```



Knowledge Bases

```
retriever = vector_db.as_retriever()  
tool = Tool(retriever, name, description)
```



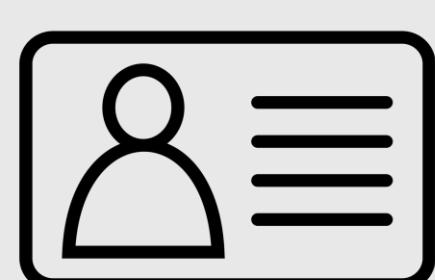
Tools – 2/3

```
def get_stock_price(ticker: str) -> float:  
    """Fetch the latest stock price for a given ticker symbol from Yahoo Finance"""  
  
    stock = yf.Ticker(ticker)  
  
    return stock.history(period="1d")["Close"].iloc[-1]  
  
    return go(f, seed, [])  
}
```

Name

Description

Content



Tools – 3/3



```
{  
    "name": "get_stock_price",  
    "description": "Fetch the stock prices  
given a ticker."  
    "parameters": {...}  
}
```

Thought: I need a tool to solve this question.
Let's see which tool can address the query.

Invoking get_stock_price...

“..according to article XX...”

I now know the final answer.



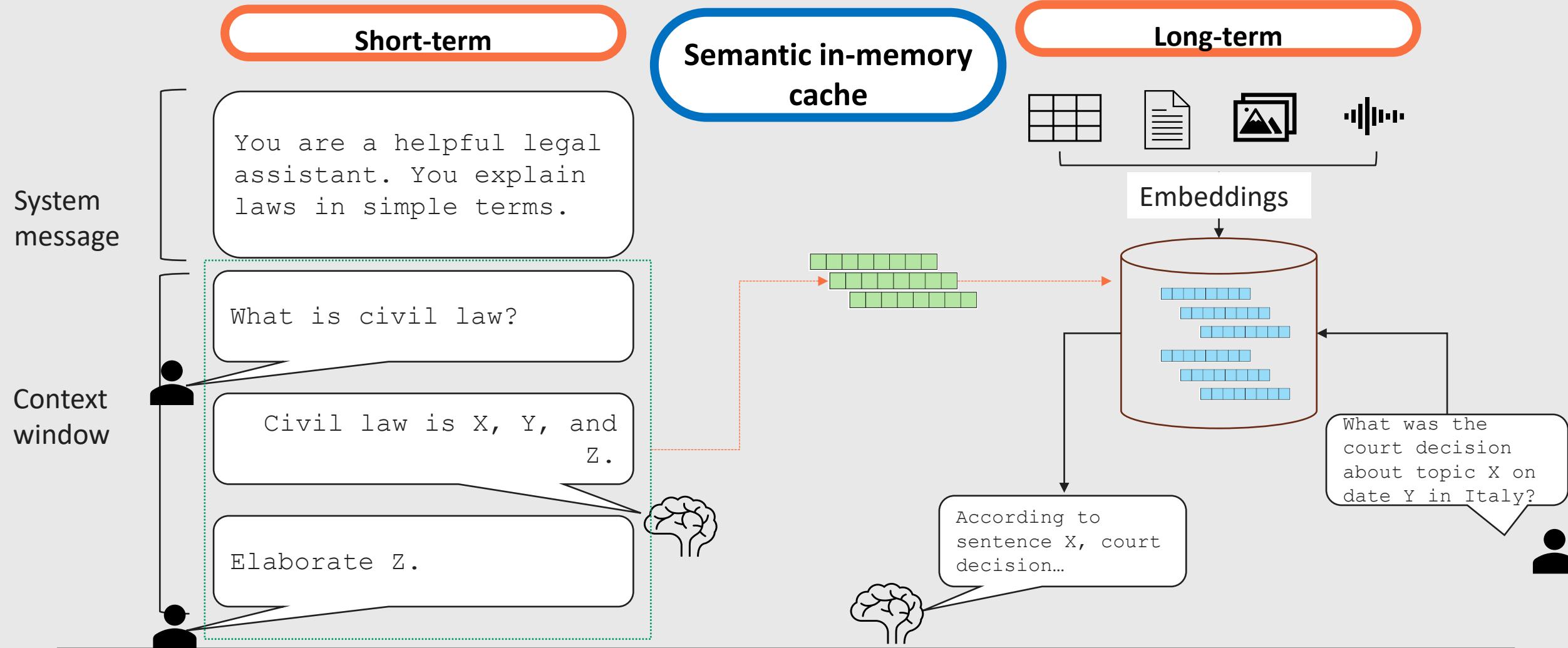
According to
sentence X, court
decision...

What's the price
of MSFT stock
today?



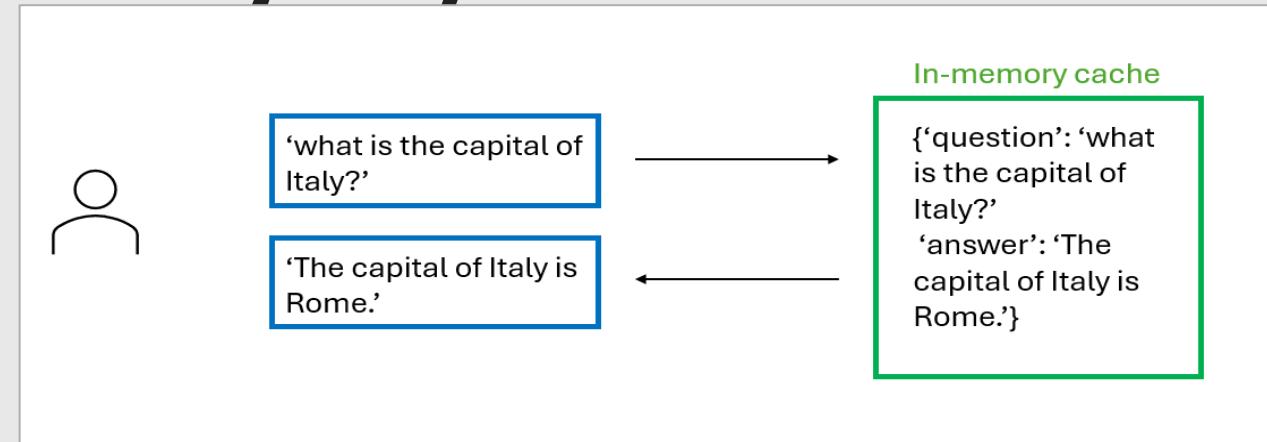
You are a helpful Financial
Assistant assistant.
You have access to the
following tools:
{tools}

Memory - 1/2

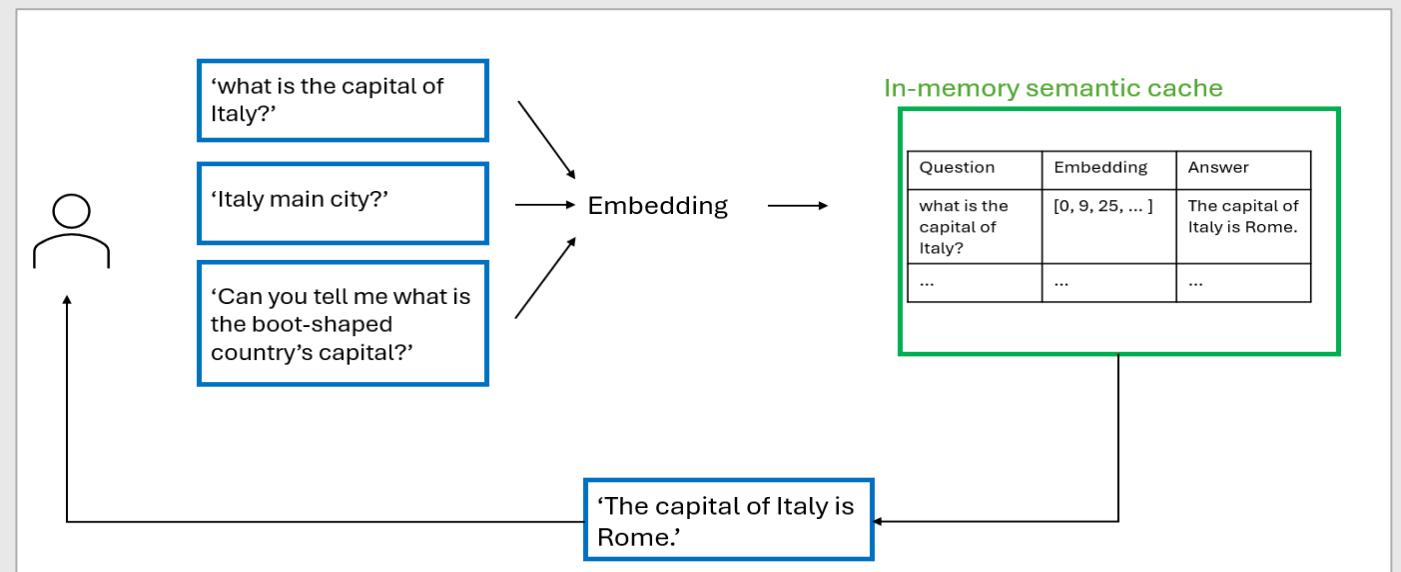


Memory – 2/2

Key-word in-memory cache



Semantic in-memory cache



Orchestrator



....and many more!

AI orchestrators provide the building blocks for designing, structuring, and managing intelligent workflows—enabling developers to compose multi-step reasoning, tool use, and agent collaboration into modular, scalable applications.

 Abstractions for LLM Workflows

 Pre-Built Integrations

 Memory and State Management

 Multi-agent Collaboration

 Debugging, Testing, and Evaluation Tools

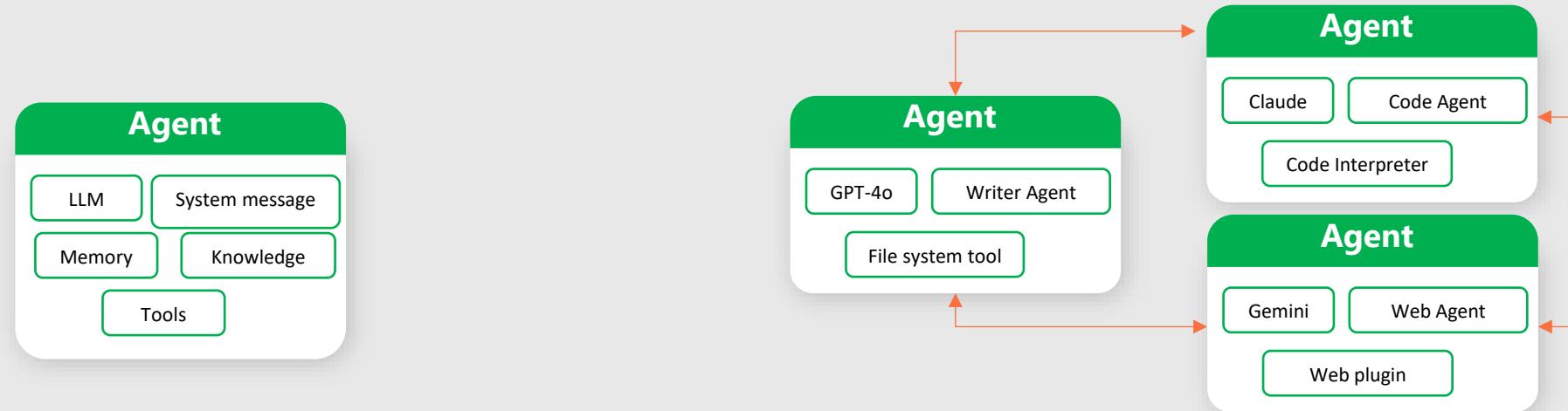
 Extensibility

03

What Happens If We Put Multiple Agents into the Same Room?



Scaling Your AI Agent



From single agent.....

Too many tools, too
hallucinations

Growing context

Hard to handle
dynamic tasks

...to multi-agent

Manageability

Predictability

Flexibility

Agents' ID Card



```
● ● ●  
web_agent = InitializeAgent(  
    name="web_agent",  
    system_message="""  
        You are an expert web surfer.  
        Given the user's query, you find the best updated content from the web.  
    """,  
    description="agent that scrapes the web and retrieves up to date information",  
    tools = [web_scraping_tool, captcha_tool]  
)
```

Name

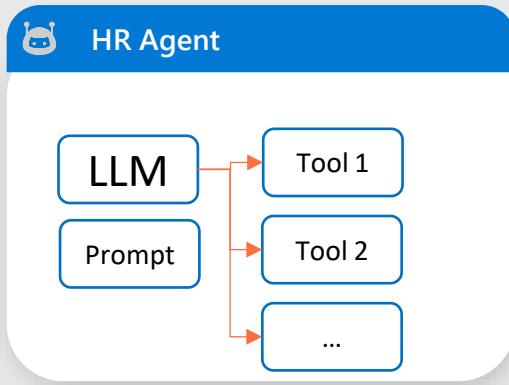
System message

Description

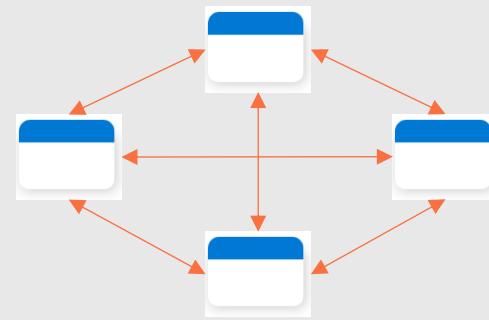
Tools

Agents Orchestration and Communication Styles

Single Agent



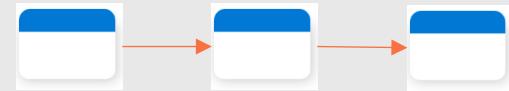
Network



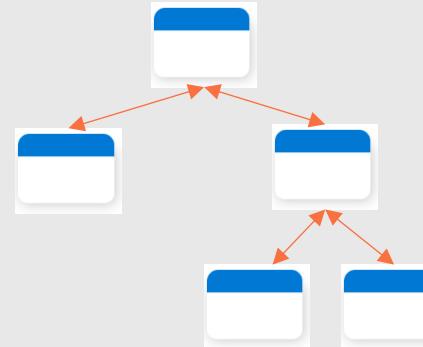
Reflection



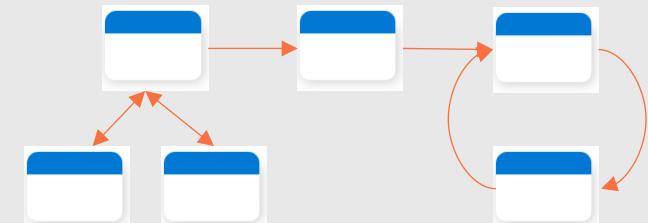
Sequential



Hierarchical



Hybrid



08

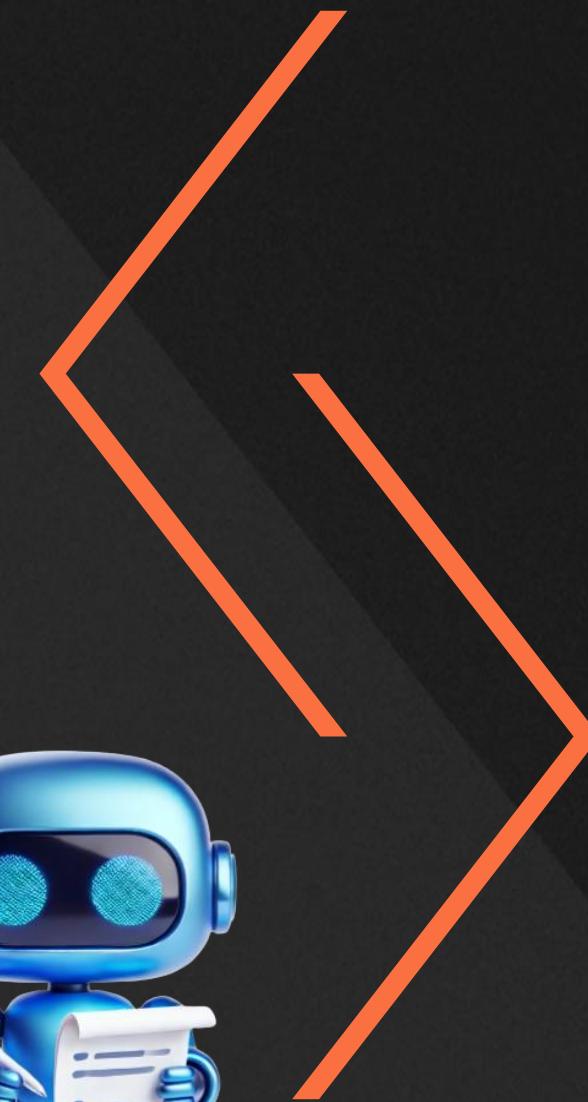
Questions



6/20/2025



X



04

Design Principles and Patterns



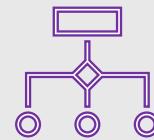
Key Principles to Keep in Mind



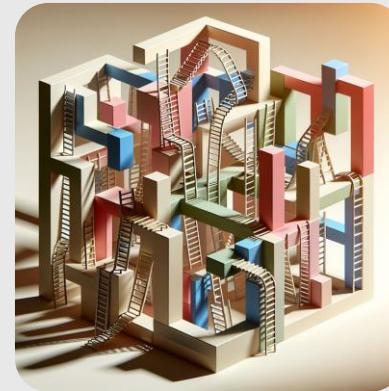
Allow them to plan and execute the best strategy



Autonomy



Breakdown and simplify complexity



Abstraction



Breakdown into smaller, reusable components.



Modularity

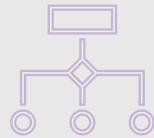
Key Principles to Keep in Mind



Allow them to plan and execute the best strategy



Autonomy



Breakdown and simplify complexity



Abstraction



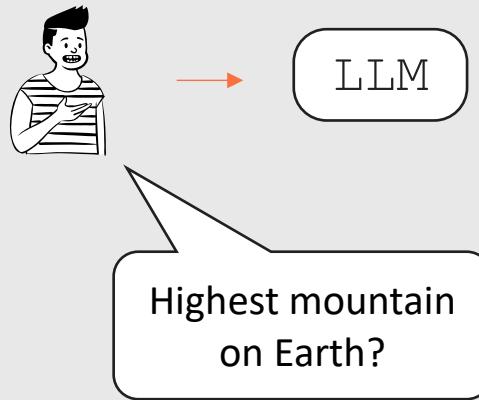
Breakdown into smaller, reusable components.



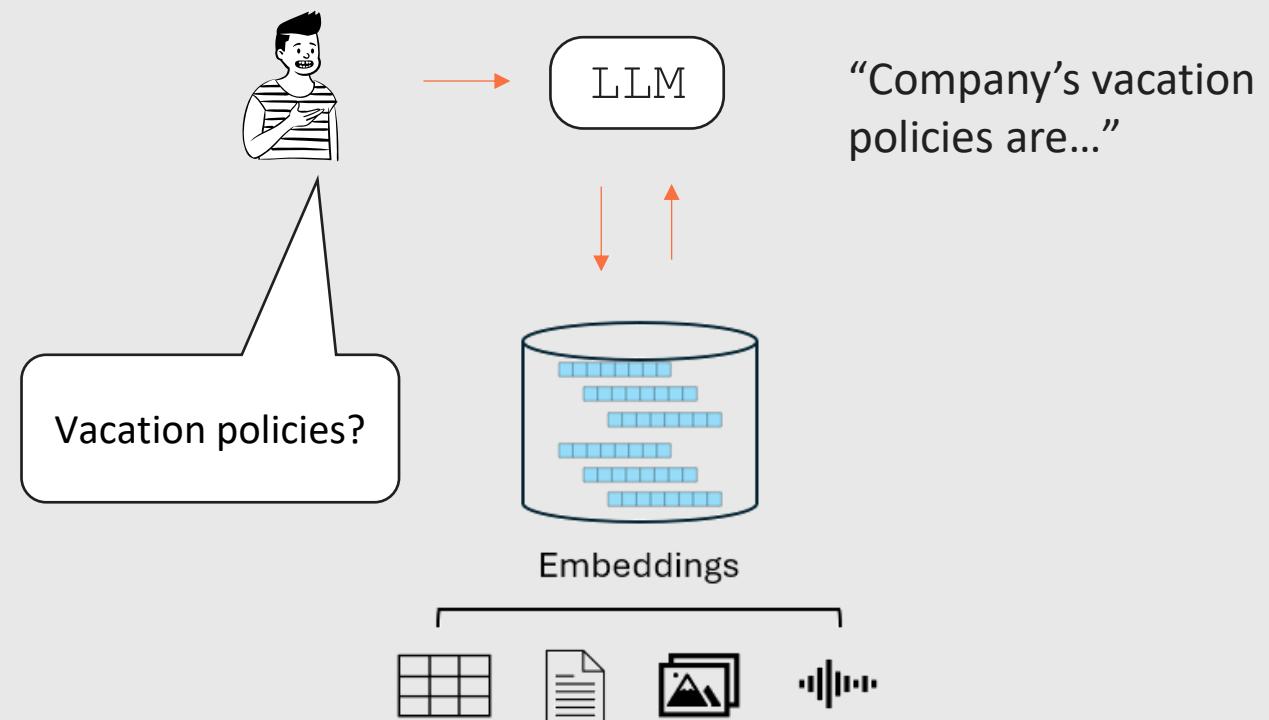
Modularity

Levels of Autonomy

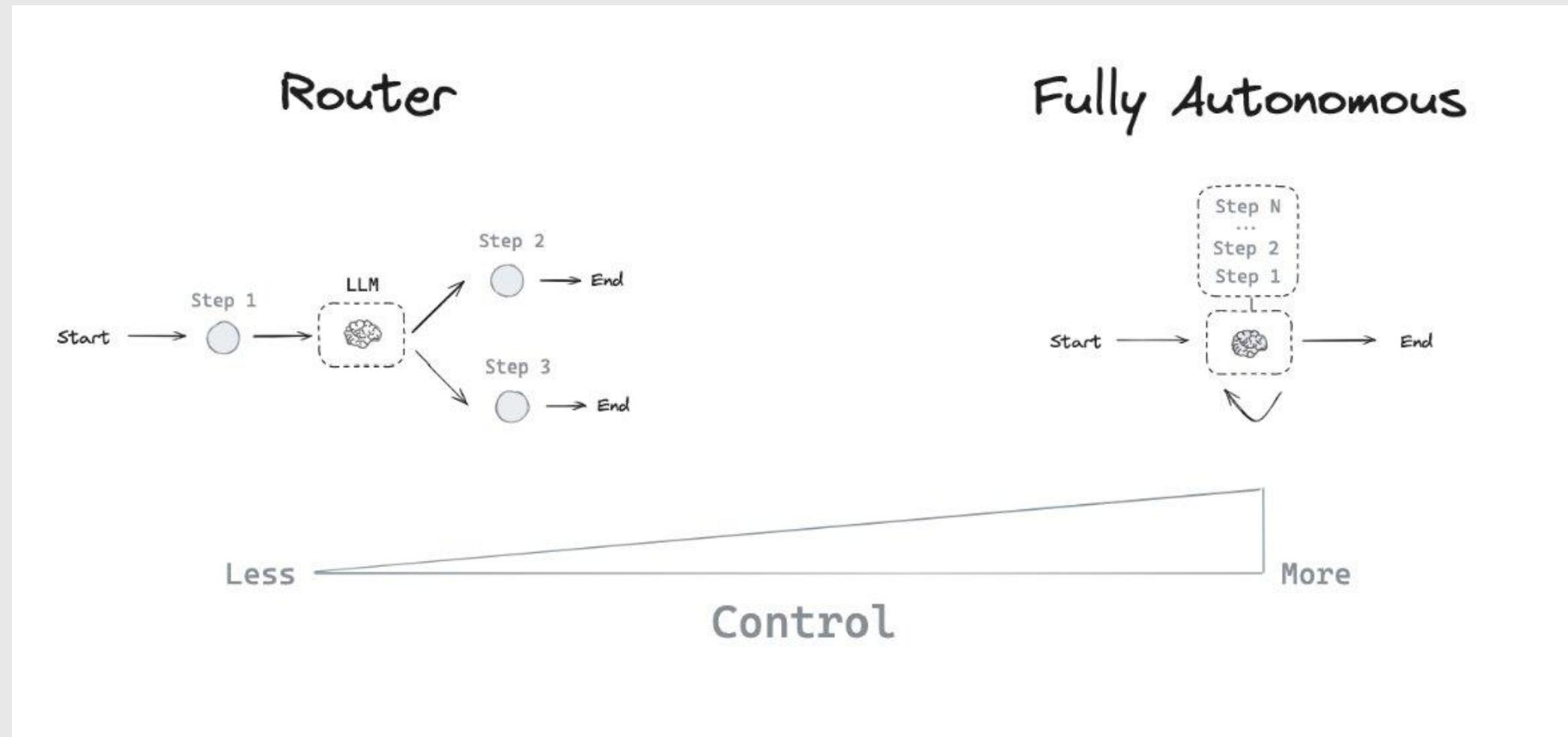
LLM API Call



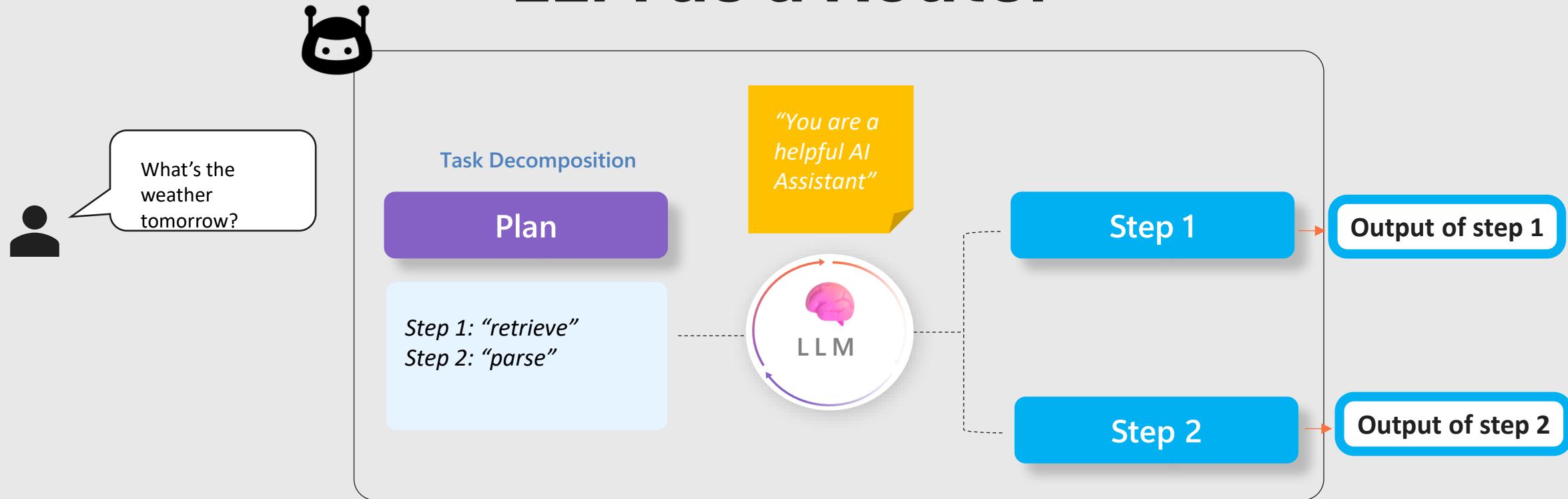
RAG



Autonomy = Degree of Freedom

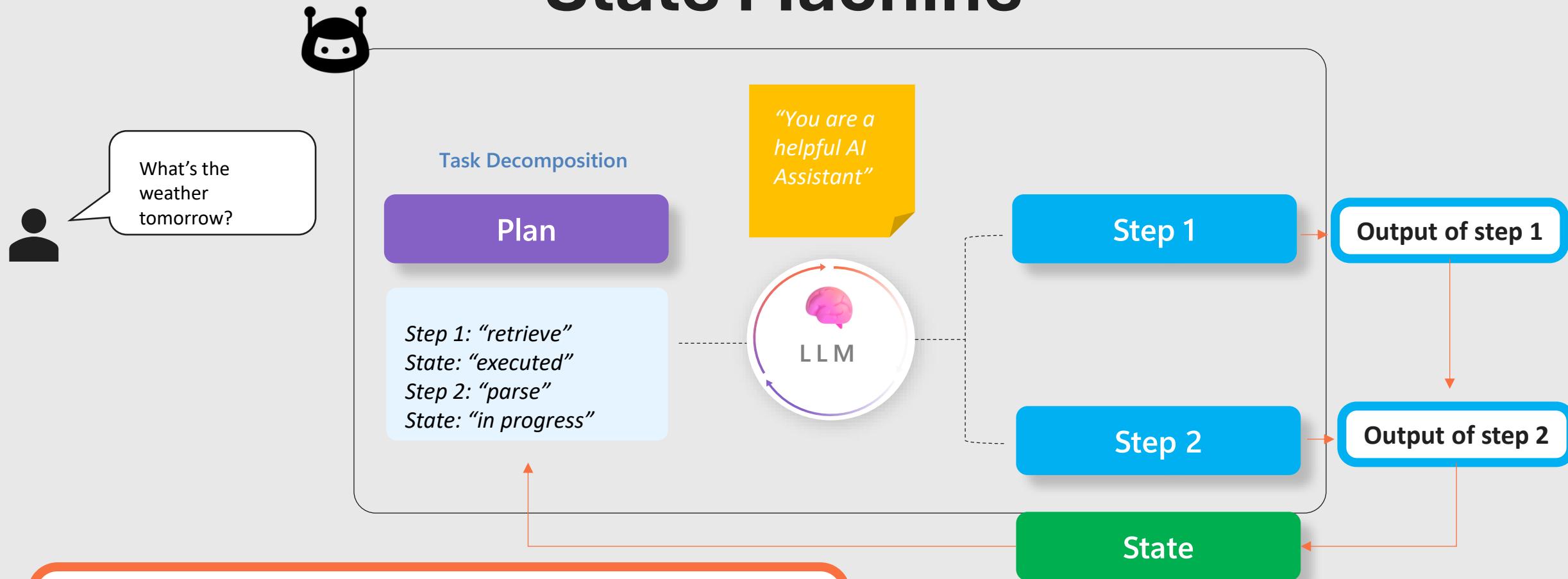


LLM as a Router



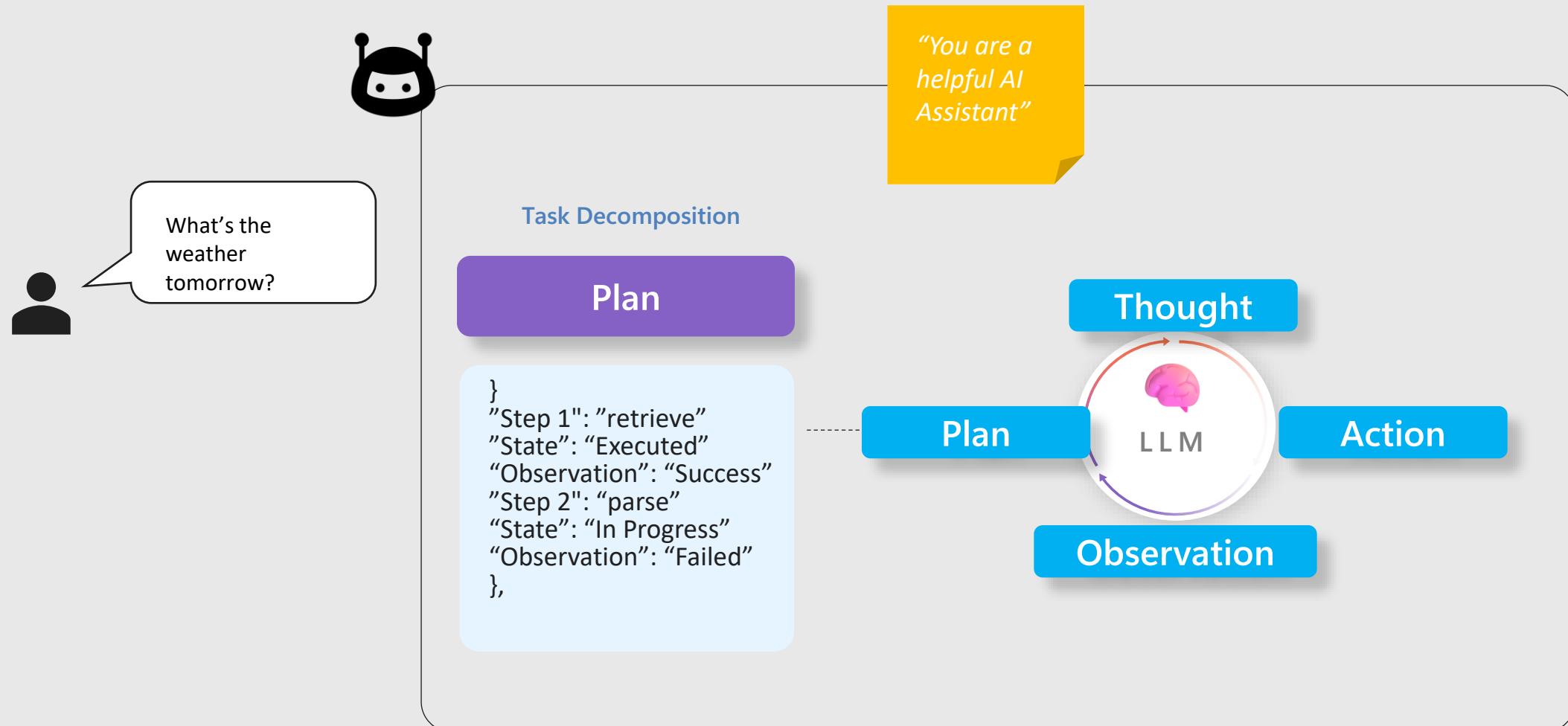
A router allows an LLM to select a single step from a specified set of options. It makes a single decision to produce a specific output. No multi-turn.

State Machine

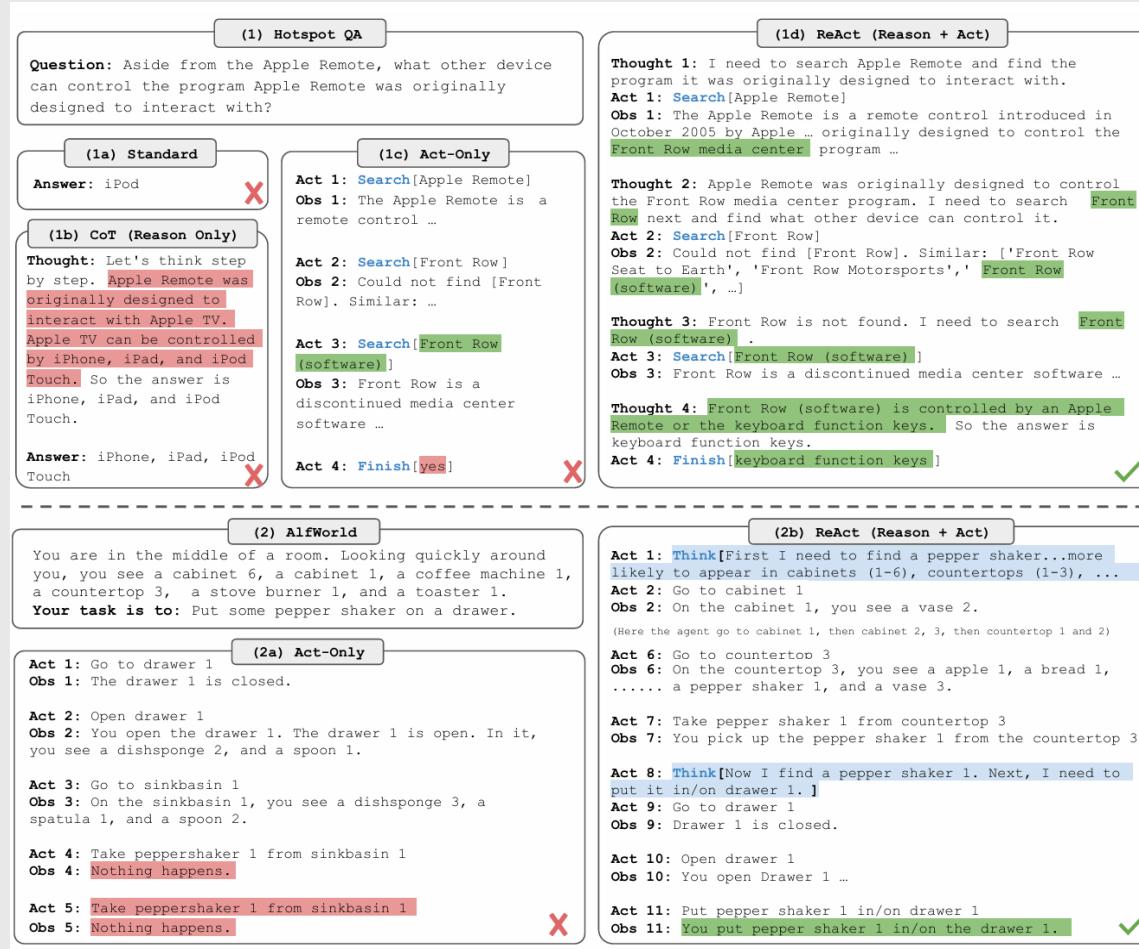


A state machine integrates an LLM with a routing loop, allowing for multi-step routing. The concept of State is part of the LangGraph ecosystem.

ReAct – Reason and Act



ReAct – Paper

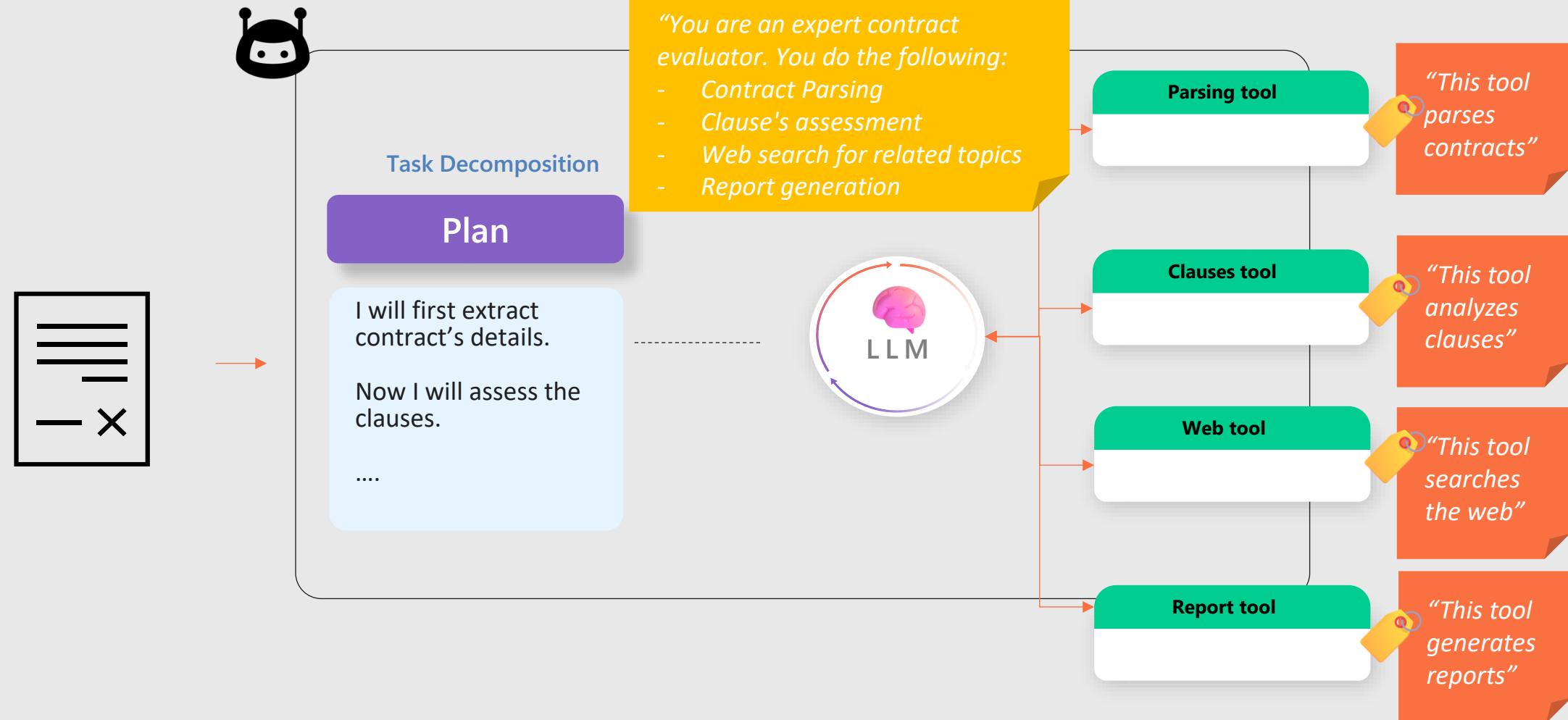


“While large language models (LLMs) have demonstrated impressive capabilities across tasks in language understanding and interactive decision making, their abilities for reasoning (e.g., chain-of-thought prompting) and acting (e.g., action plan generation) have primarily been studied as separate topics.

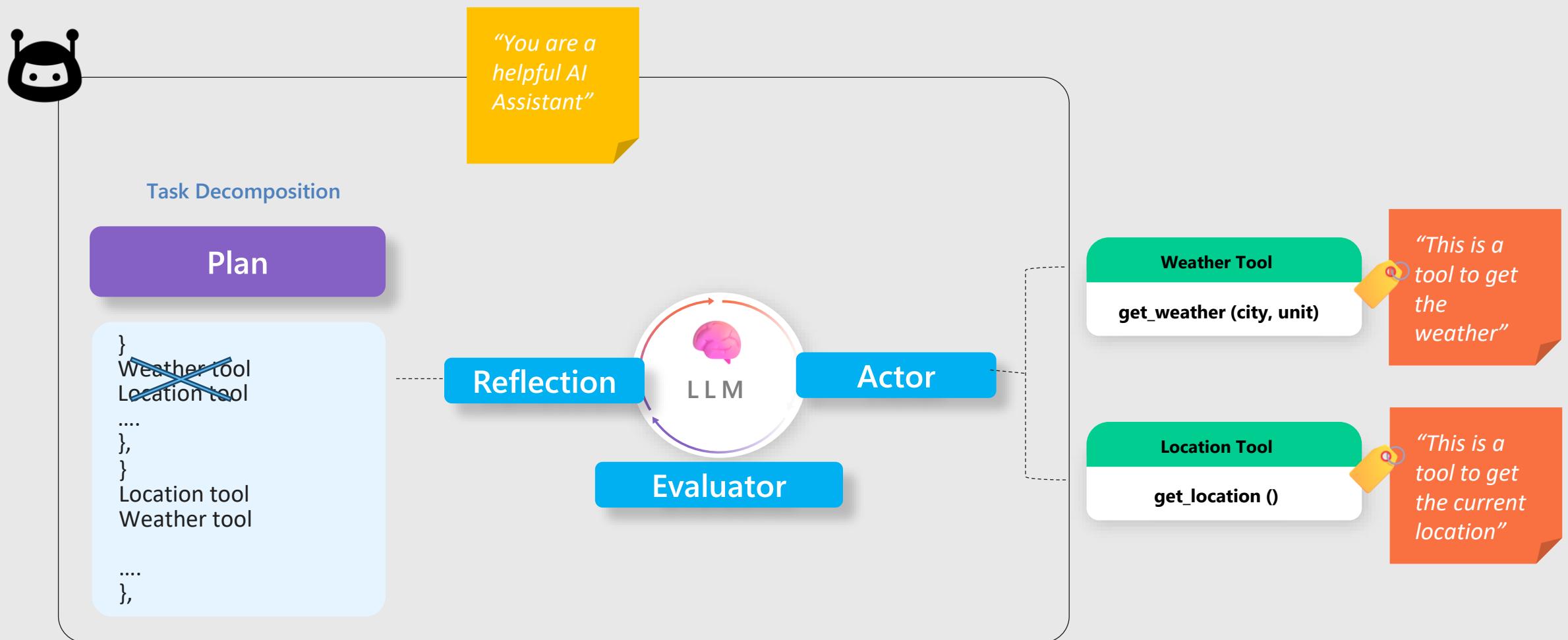
[...] reasoning traces help the model **induce, track, and update action plans** as well as handle exceptions, while actions allow it to **interface with external sources**, such as knowledge bases or environments, to gather additional information.”

Source: <https://arxiv.org/abs/2210.03629>

Structured Planning



Self-Reflection

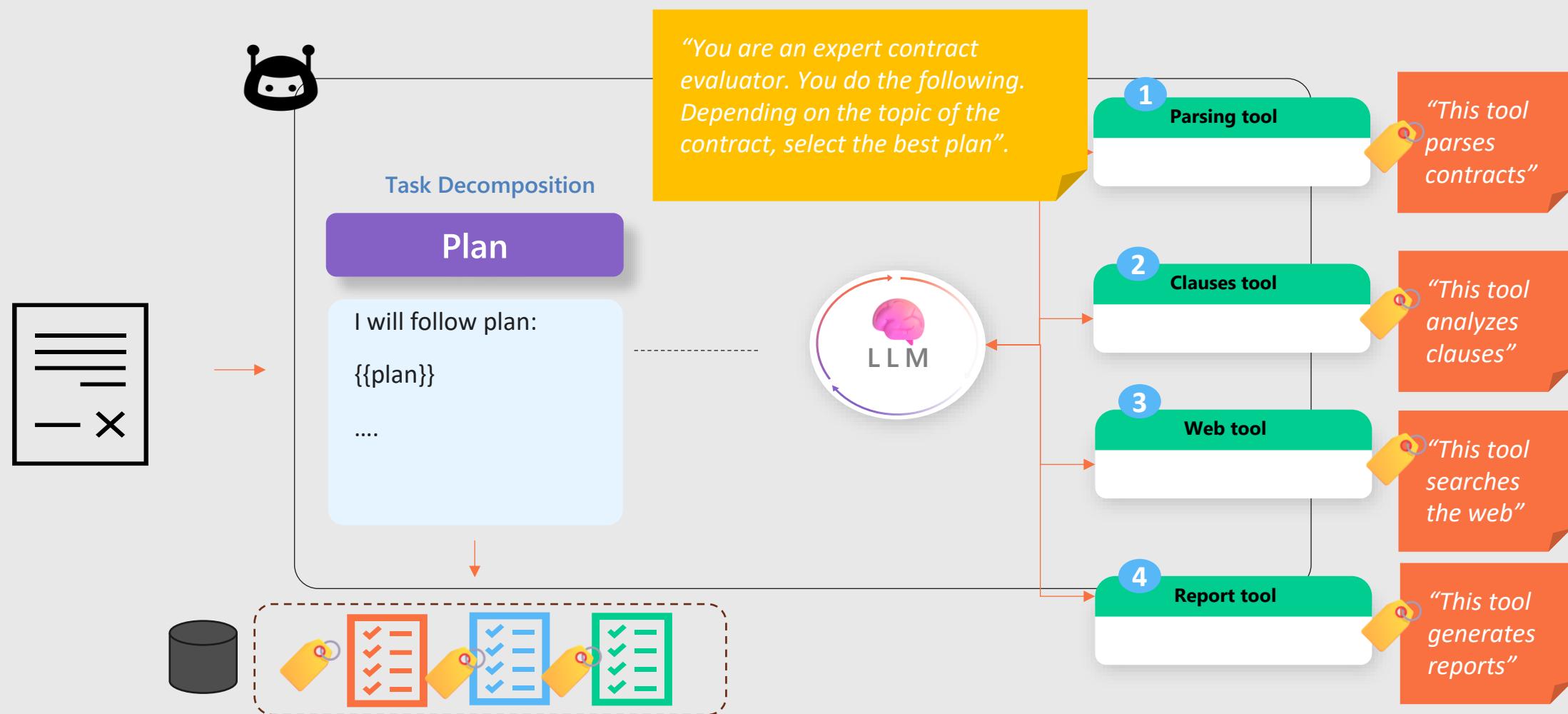


Source: <https://arxiv.org/pdf/2303.11366>

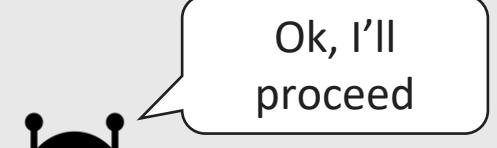
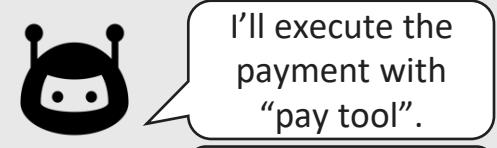
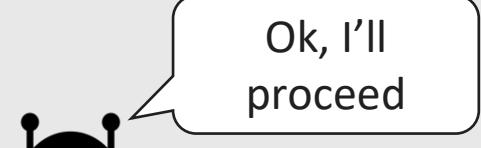
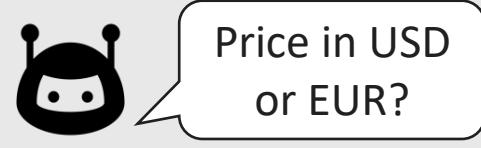
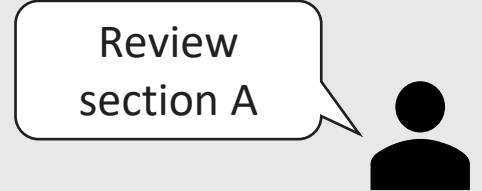
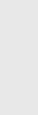
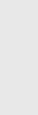
Self-Reflection – Paper

	1. Decision making	2. Programming	3. Reasoning
(a) Task	You are in the middle of a room [...] Task: clean some pan and put it in countertop.	Task: You are given a list of two strings [...] of open '(' or close ')' parentheses only [...]	Task: What profession does John Lanchester and Alan Dean Foster have in common?
(b) Trajectory	[...] Action: take pan1 from stoveburner1 Obs: Nothing happens. [...] Action: clean pan1 with sinkbasin1 Obs: Nothing happens. [...]	def match_parens(lst): if s1.count('(') + s2.count('(') == s1.count(')') + s2.count(')'): [...] return 'No'	Think: [...] novelist, journalist, critic [...] novelist, screenwriter [...] common is novelist and screenwriter. Action: "novelist, screenwriter"
(c) Evaluation (internal / external)	Rule/LM Heuristic: Hallucination.	Self-generated unit tests fail: assert match_parens(...)	Environment Binary Reward: 0
(d) Reflection	[...] tried to pick up the pan in stoveburner 1 [...] but the pan was not in stoveburner 1. [...]	[...] wrong because it only checks if the total count of open and close parentheses is equal [...] order of the parentheses [...]	[...] failed because I incorrectly assumed that they both had the same multiple professions [...] accurately identifying their professions.
(e) Next Trajectory	[...] Action: take pan 1 from stoveburner 2 [...] Obs: You put the pan 1 in countertop 1.	[...] return 'Yes' if check(S1) or check(S2) else 'No'	Think: [...] So the profession John Lanchester and Alan Dean Foster have in common is novelist. Action: "novelist"

Plan Selection



Human in the Loop



Validation station

Feedback and editing

Input required

Tool approval

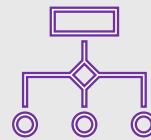
Key Principles to Keep in Mind



Allow them to plan and execute the best strategy



Autonomy



Breakdown and simplify complexity



Abstraction



Breakdown into smaller, reusable components.



Modularity

Agent Abstraction – First-class Citizen

Single Agent abstraction

LLM (GPT-4o, Llama, Claude...)

Tools

System message (“persona”)

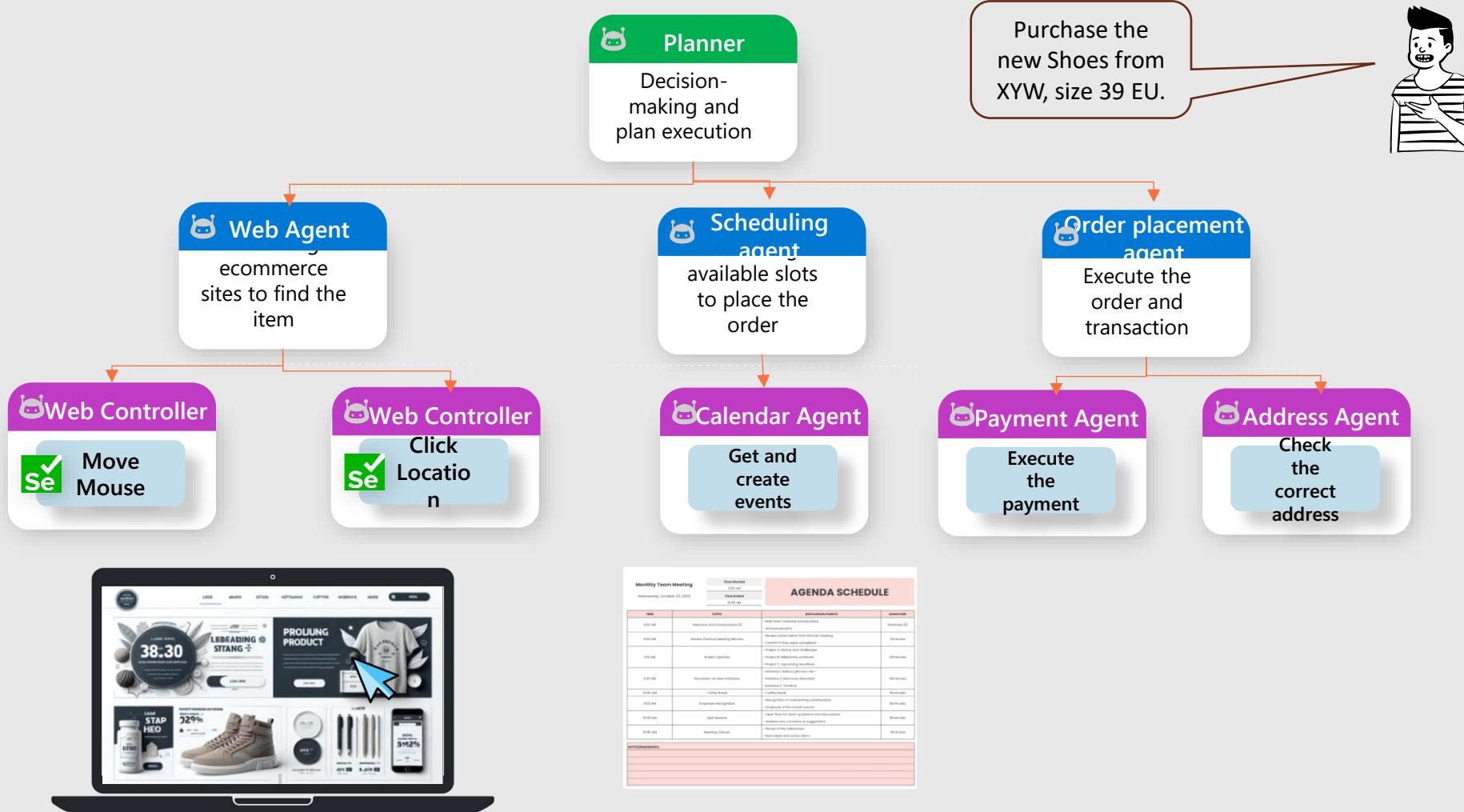
Multi-agent abstraction

Many agents working together at different level of management

Different layers of execution

Different degrees of autonomy

Modularity and Abstraction



Operator

Find me a family friendly campsite at Joshu|

0

Key Principles to Keep in Mind



Allow them to plan and execute the best strategy



Autonomy



Breakdown and simplify complexity



Abstraction

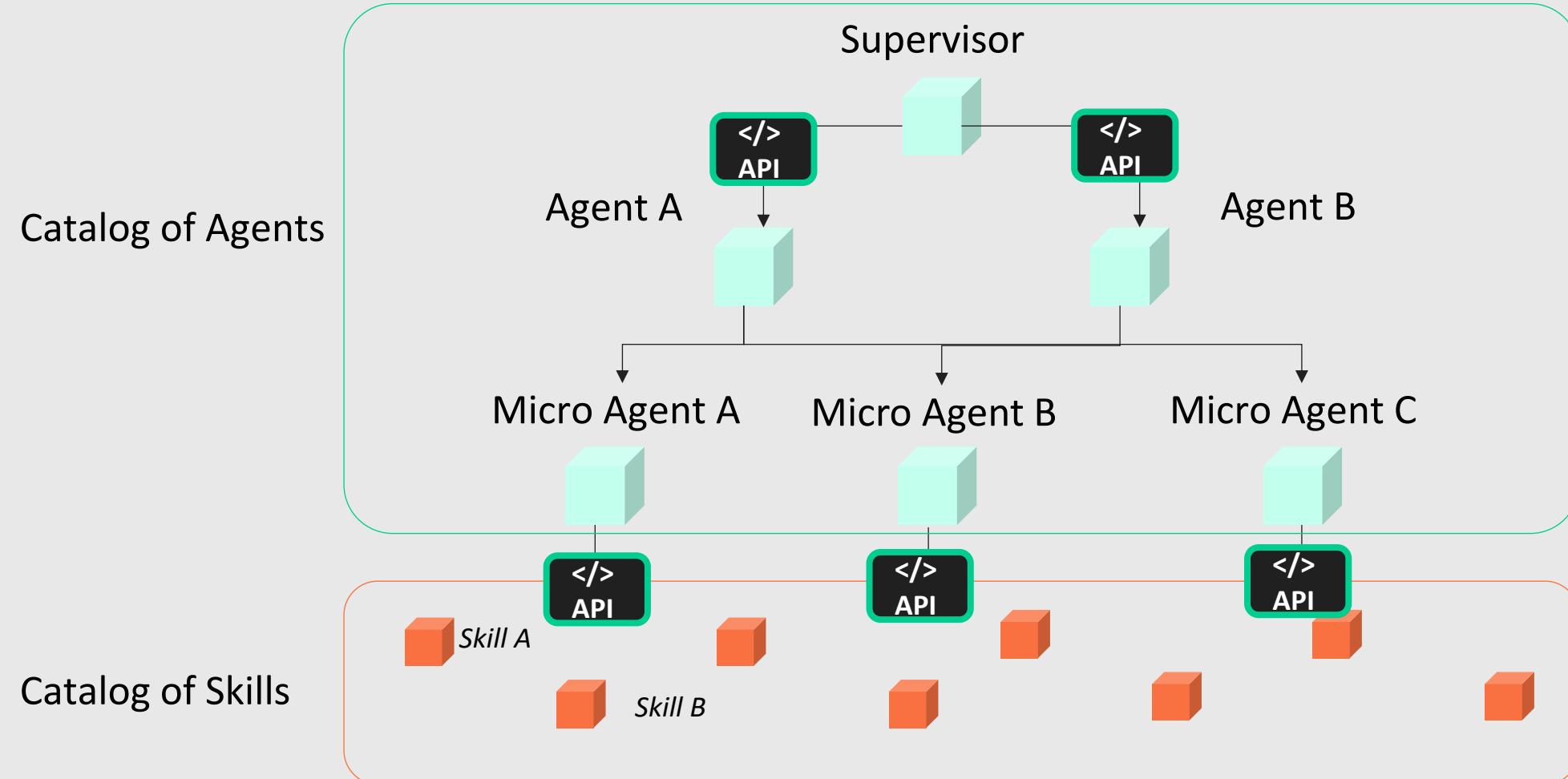


Breakdown into smaller, reusable components.



Modularity

Microservices Go Hand in Hand with Modularity and Abstraction



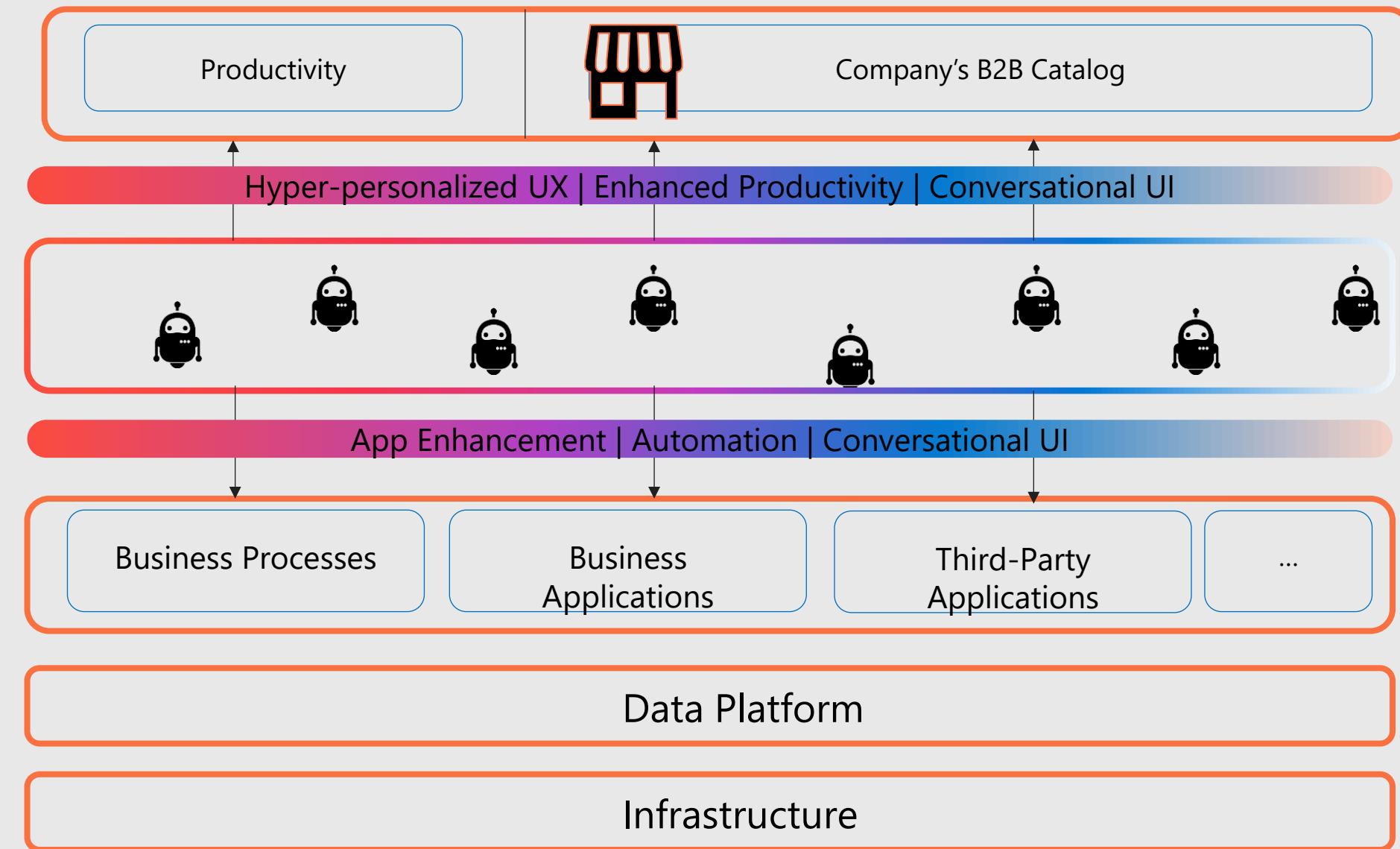
Engagement Platform (internal and external)

Agent Orchestration

Agentic State

Agent Orchestration

Enterprise Applications



08

Questions



6/20/2025



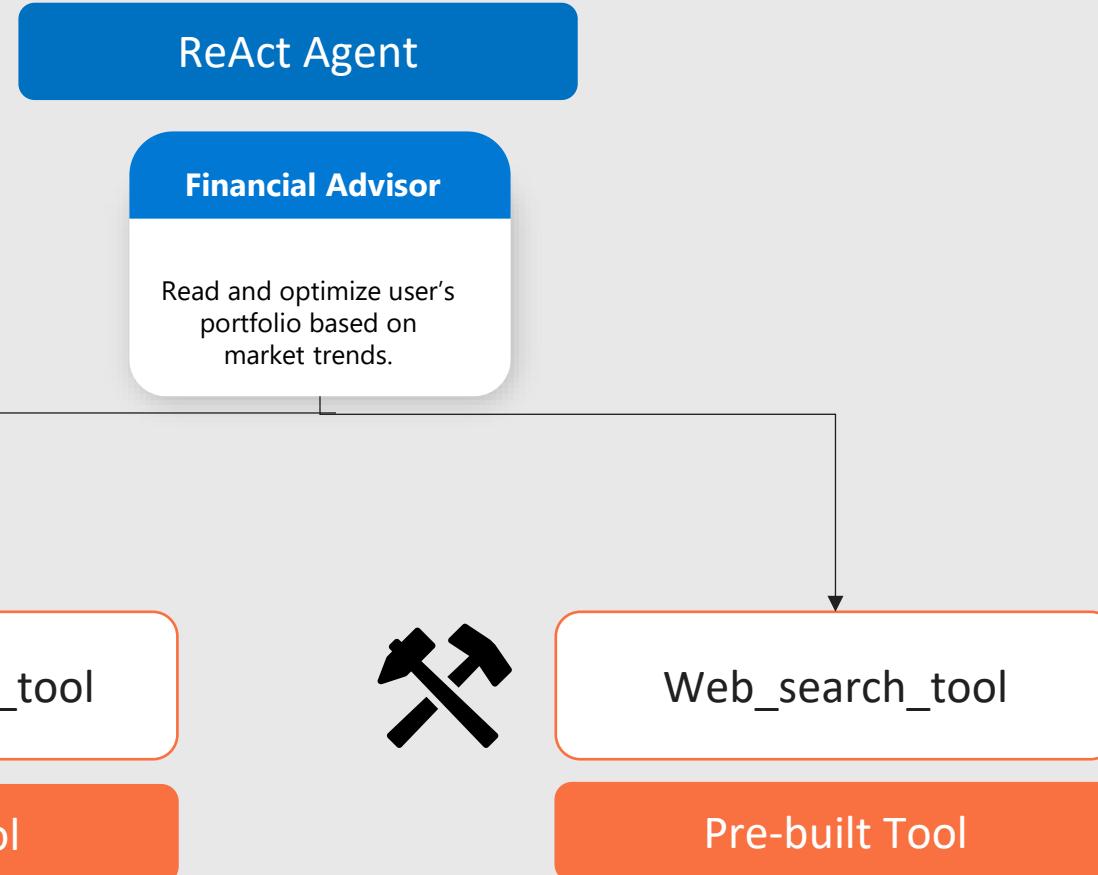
X



07

Concept Demo: Build your first "mini agent"







x

AI Agents Open Protocols

AI Agents Open Protocols

**Model Context
Protocol**



Agent to Agent



NLWeb



AI Agents Open Protocols

**Model Context
Protocol**



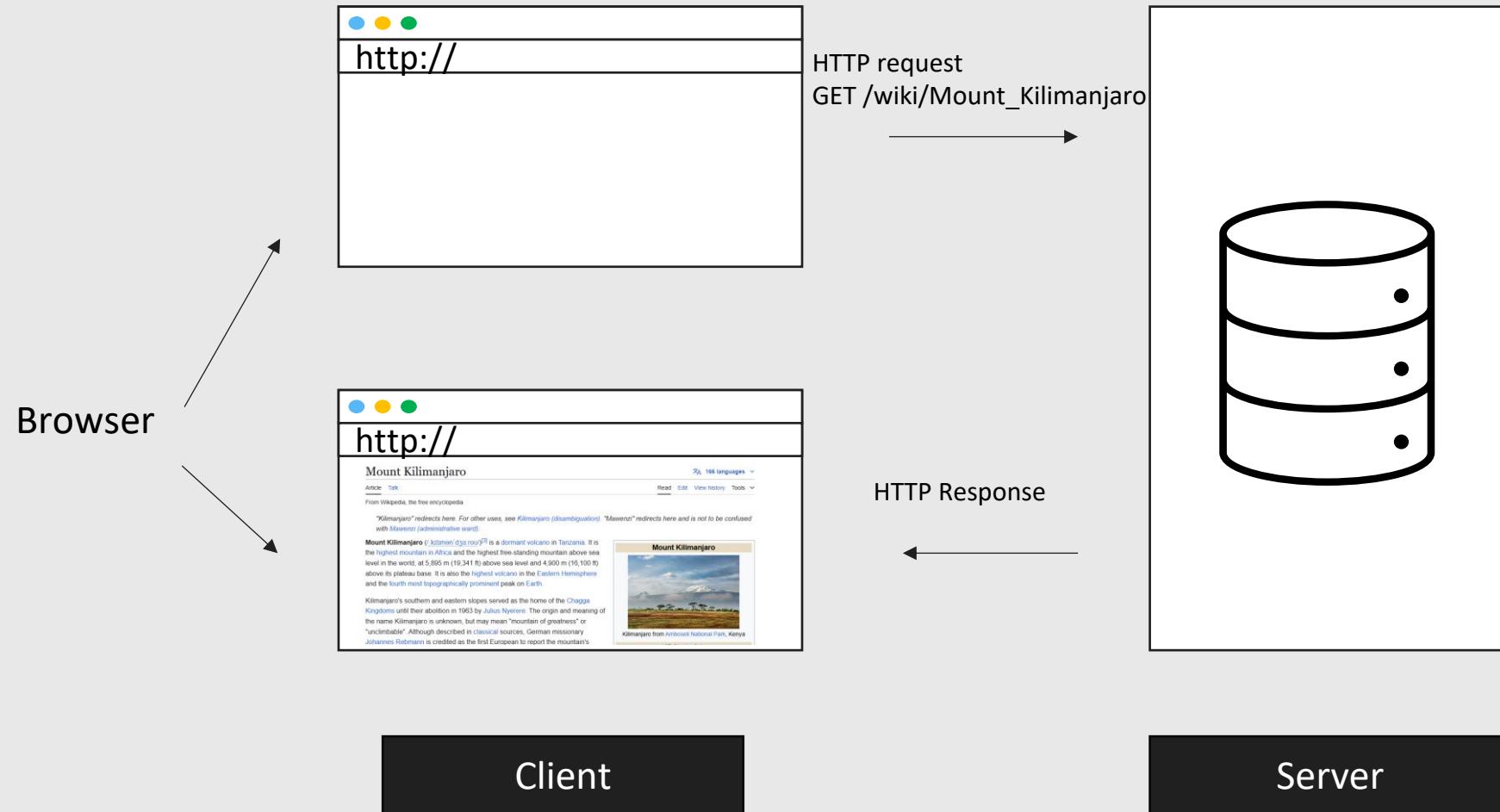
Agent to Agent



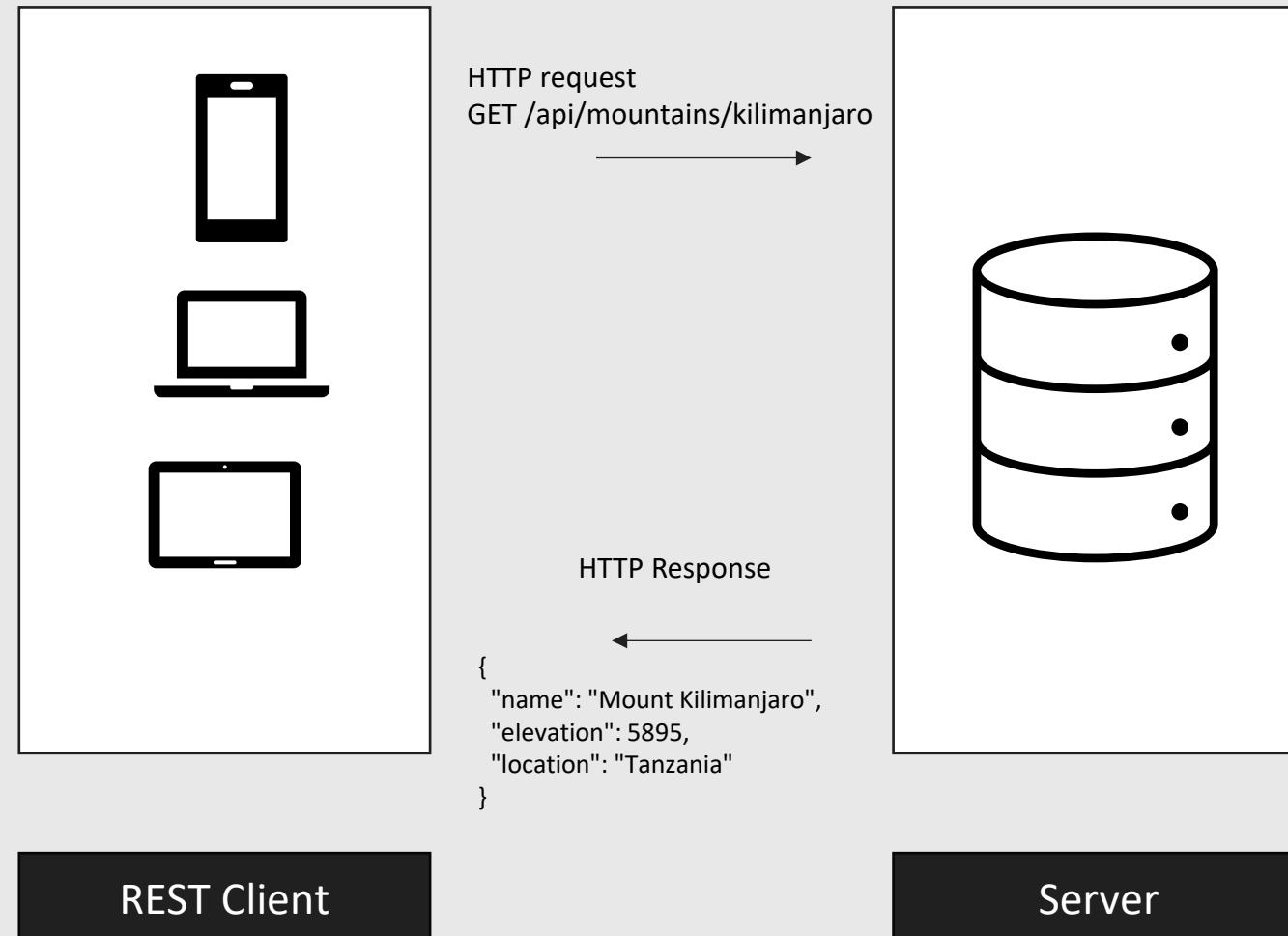
NLWeb



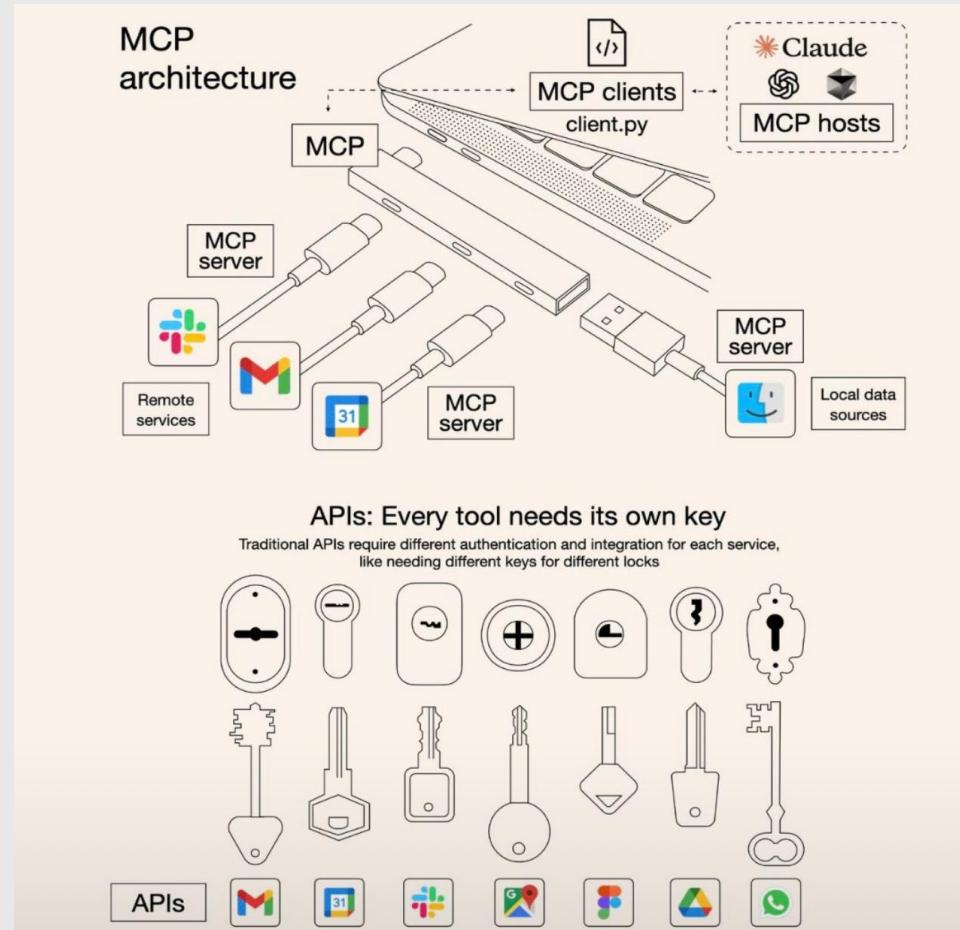
What is a protocol?



What is an API?



Model Context Protocol

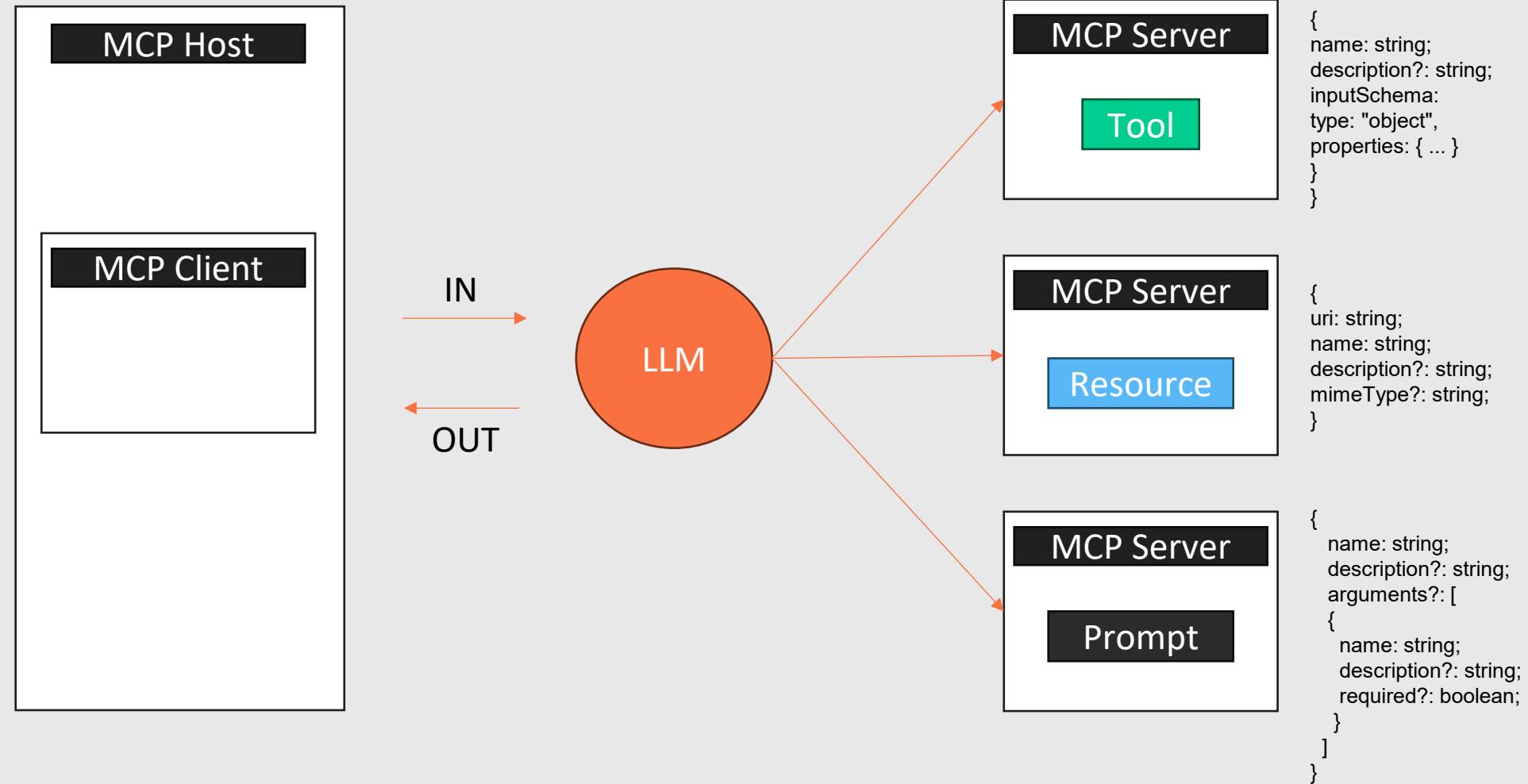


Source: <https://blog.csdn.net/ChailangCompany/article/details/146354913>

"MCP is an open protocol that standardizes how applications provide context to LLMs. Think of MCP like a USB-C port for AI applications. Just as USB-C provides a standardized way to connect your devices to various peripherals and accessories, MCP provides a standardized way to connect AI models to different data sources and tools."

Source: <https://modelcontextprotocol.io/introduction>

High-Level Architecture of MCP



Example – Weather Forecast MCP Server

```
● ● ●

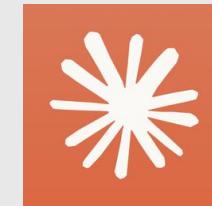
# server.py
from mcp.server.fastmcp import FastMCP
import yfinance as yf

# Create an MCP server
mcp = FastMCP("Demo")

# Add an addition tool
@mcp.tool()

def get_stock_price(ticker: str) -> float:
    """Fetch the latest stock price for a given ticker symbol from Yahoo Finance"""
    stock = yf.Ticker(ticker)
    return stock.history(period="1d")["Close"].iloc[-1]

if __name__ == "__main__":
    # Initialize and run the server
    mcp.run(transport='stdio')
    return go(f, seed, [])
}
```



Claude Desktop
as MCP Host



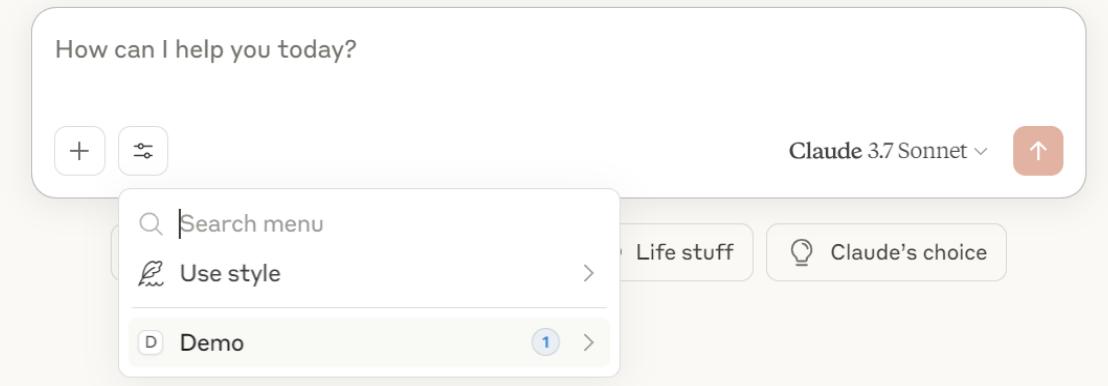
Local Machine as
MCP Server,
using Python
SDK

Tip: Find a catalog of MCP Servers here!
[Source: https://mcpservers.org/](https://mcpservers.org/)

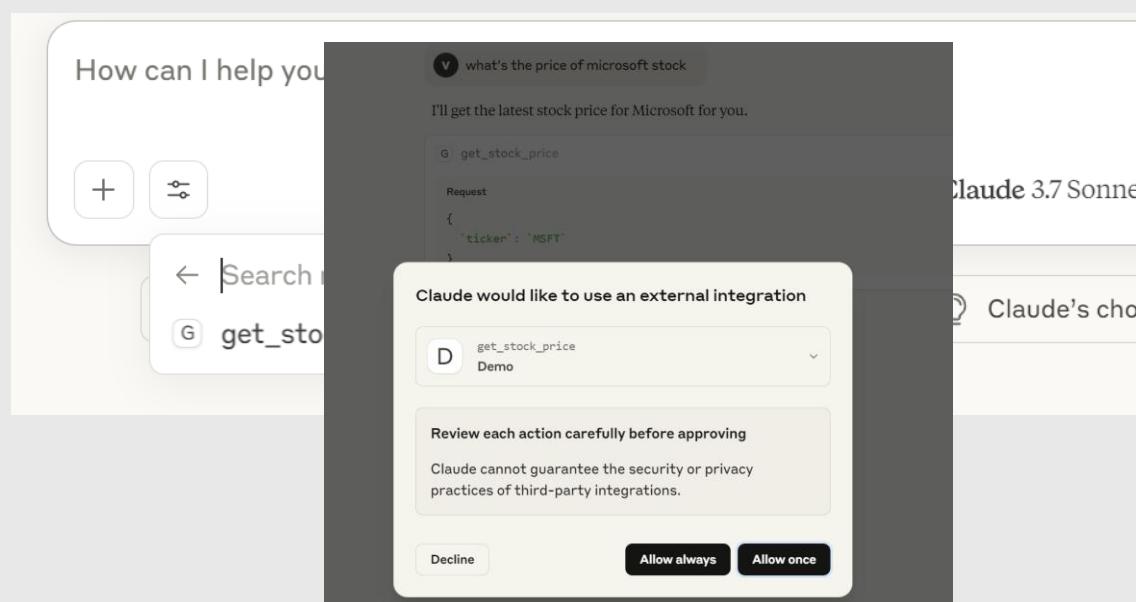


* Coffee and Claude time?

How can I help you today?



How can I help you?



v what's the price of microsoft stock

I'll get the latest stock price for Microsoft for you.

G get_stock_price

Request

```
{
  `ticker`: `MSFT`
}
```

Response

453.1300048828125

The current price of Microsoft (MSFT) stock is \$453.13 as of today, May 16, 2025.

Retry ▾

Claude can make mistakes. Please double-check responses.

AI Agents Open Protocols

Model Context
Protocol



Agent to Agent



NLWeb



Agent to Agent

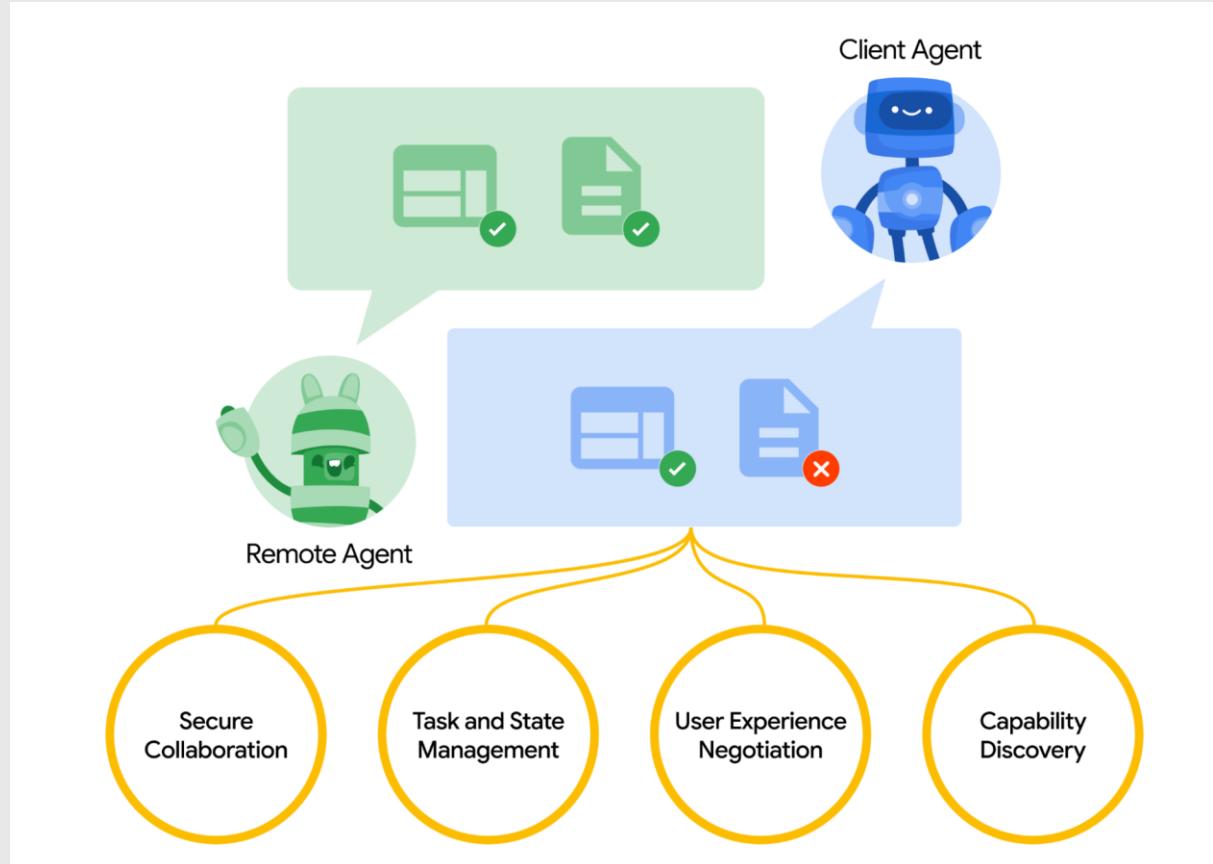


Image source: <https://a2aproto.col.ai/>

"A2A Protocol is an open standard that enables AI agents to communicate and collaborate across different platforms and frameworks, regardless of their underlying technologies. It's designed to maximize the benefits of agentic AI by enabling true multi-agent scenarios."

Source: <https://a2aproto.col.ai/>

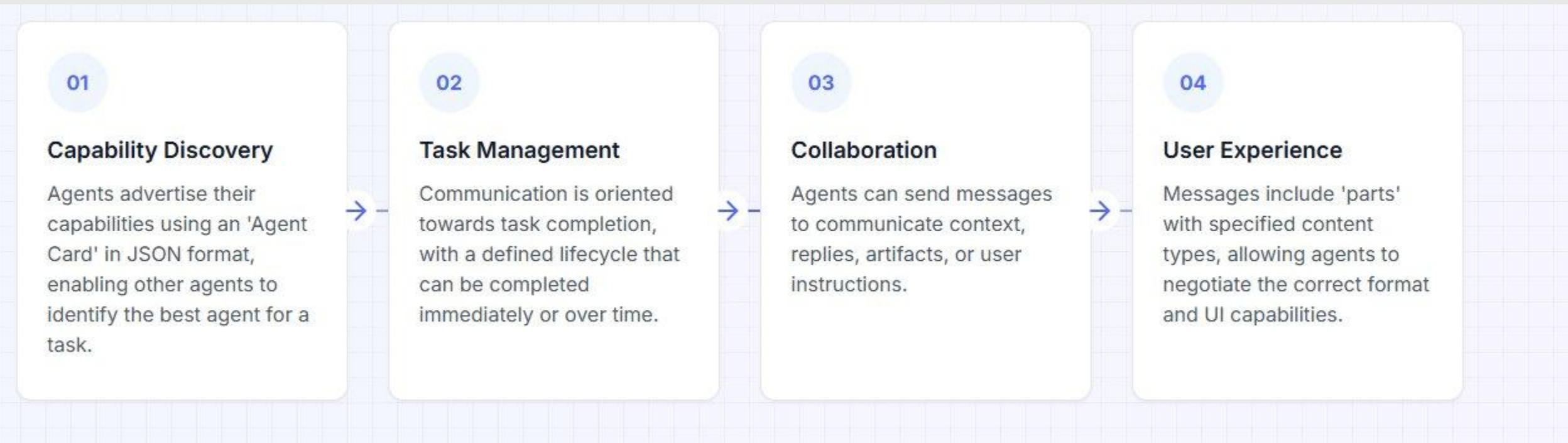
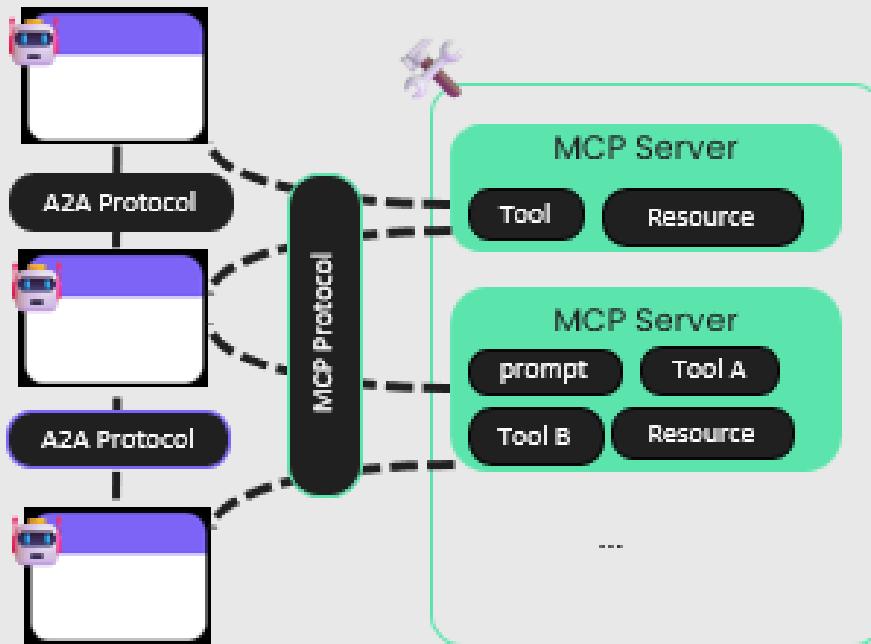


Image source: <https://a2aproto.col.ai/>

A2A vs MCP



MCP

Protocol to connect LLMs with Data, Resources and Tools.

A2A

Protocol to connect and enable collaboration among multiple agents.

AI Agents Open Protocols

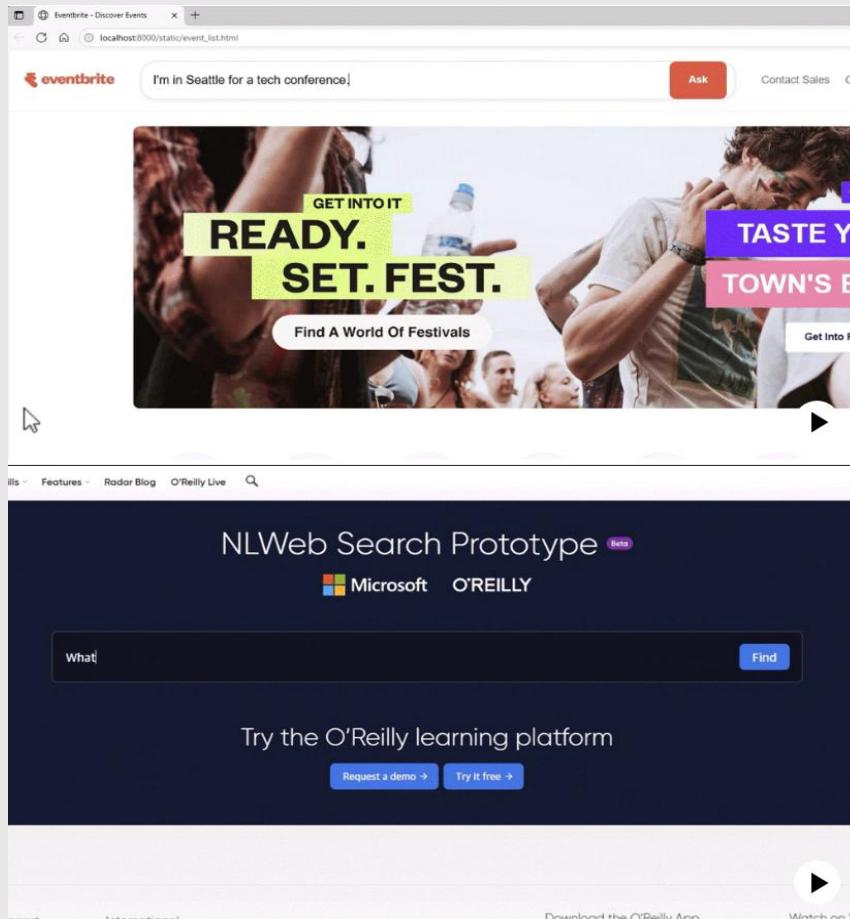
Model Context
Protocol

Agent to Agent

NLWeb



What's NLWeb?

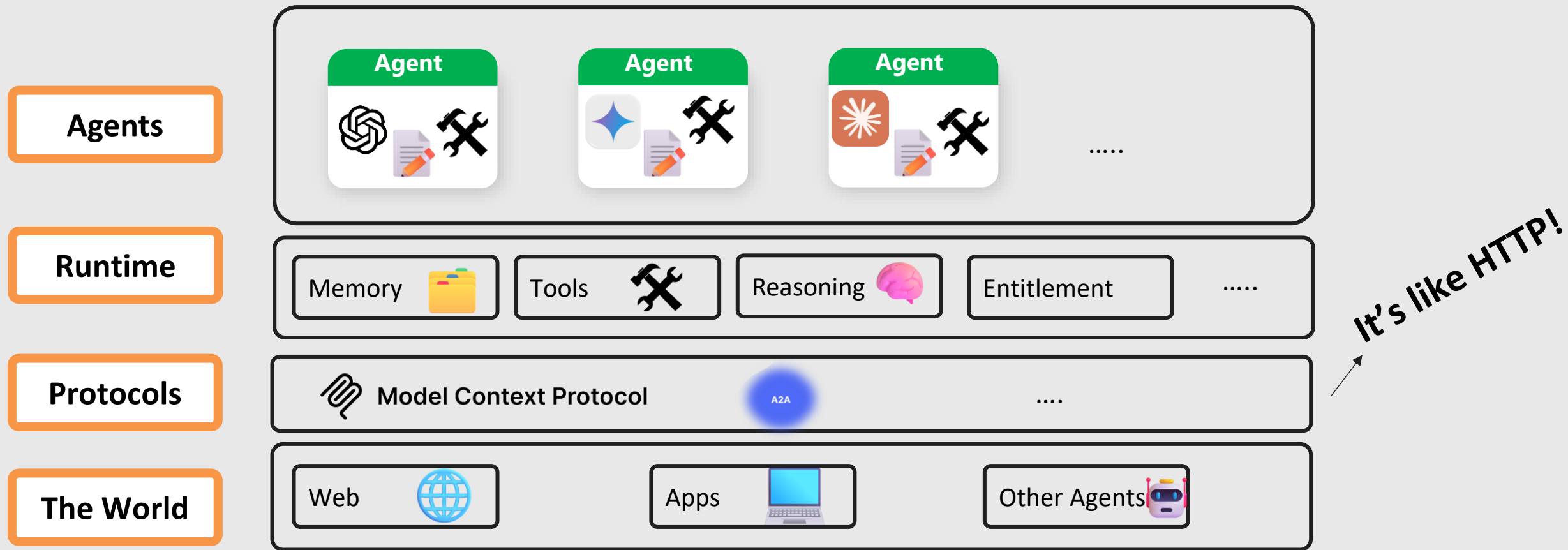


Source: <https://news.microsoft.com/source/features/company-news/introducing-nlweb-bringing-conversational-interfaces-directly-to-the-web/?msocid=16ccd88e67d963500e5fc7866d062ac>

“Building conversational interfaces for websites is hard. NLWeb seeks to make it easy for websites to do this. And since NLWeb natively speaks MCP, the same natural language APIs can be used both by humans and agents.”

Source: <https://github.com/microsoft/nlweb>

Building the Open Agentic Web



Source: <https://github.com/microsoft/nlweb>

08

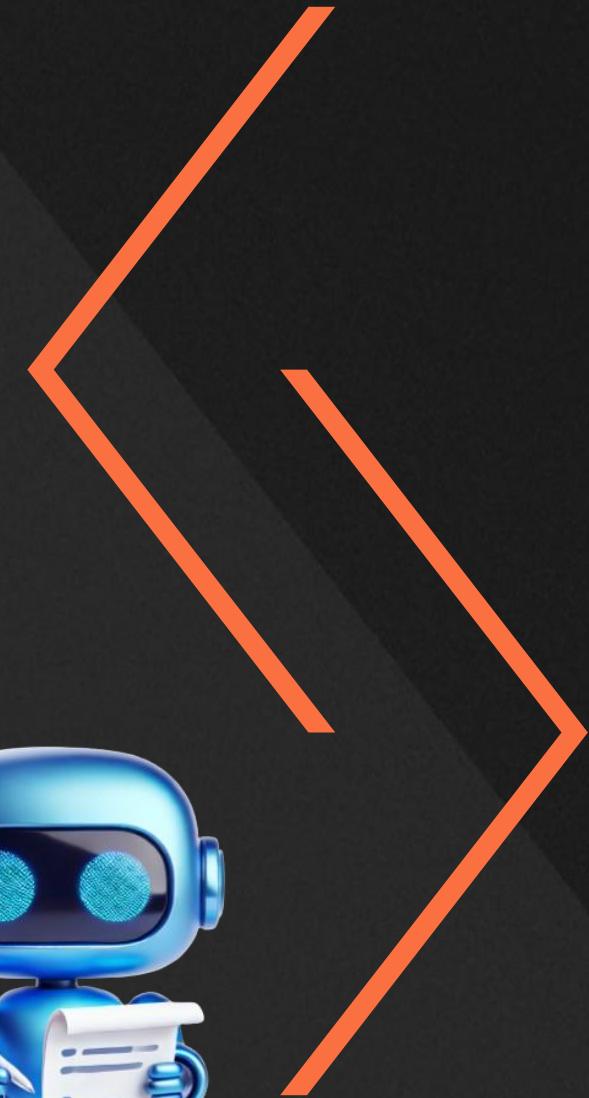
Questions



6/20/2025



X





09

Thank You!

