

Software &
Digital Platforms

AI Agent Workshop

Building AI Agents with Azure

May 05 – May 12

Contents

Introduction to Agentic AI

Day 1: May 5, 2025

01

Introduction to Agentic AI

02

Agentic System Design Patterns

03

Azure Agentic Services & Frameworks

- Azure AI Agent Service
- Sematic Kernel
- Autogen

Workshop Assignments

Day 2: May 7, 2025

04

Guided exercises to build and test agents

05

Workshop Challenges Overview

06

Team Assignments and handover for independent work

Workshop Results

Day 3: May 12, 2025

07

Presentation of Team Solutions

08

Solution Review

09

Wrap-up & Final Discussion



Workshop Overview

By the end of the workshop you will learn

How to build AI Agents and explore the Agent [tools](#)

What are the AI Agent Platforms and how to use them

Hands-on experience building single and multi-agent systems

Getting Started

[README](#) - <https://github.com/microsoft/OpenAIWorkshop/blob/int-agentic/README.md>

Target Audience

Day 1: Executives, Product Management, Customer Success, Developers

Day 2: Developers, Product Management, Customer Success

Day 3: Developers



Pre-requisites

Day 1: Basic understanding of AI and LLMs

Day 2 & Day 3: Access to an Azure subscription with \$50 budget, [GitHub](#) account, VSCode, basic familiarity with python

Please complete before next meeting if you haven't already:

Questionnaire

[Microsoft Forms](#)

Agents quickstart

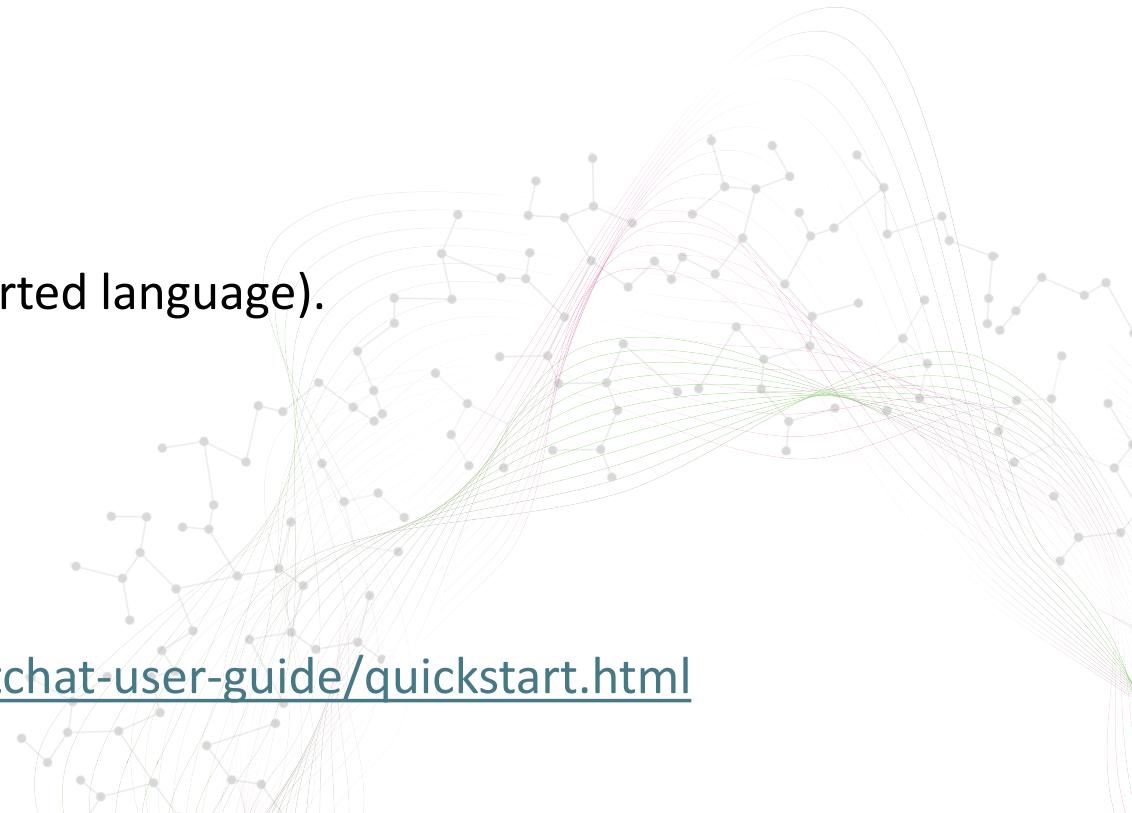
[Azure Agents Quickstart Guide](#) (Python or any other supported language).

Semantic Kernel quickstart

[Semantic Kernel Quickstart Guide](#)

AutoGen quickstart

<https://microsoft.github.io/autogen/dev/user-guide/agentchat-user-guide/quickstart.html>



Workshop Support Team



James Nguyen



Anil Dwarkanath



Nicole Serafino



Patrick O'Malley



Heena Ugale



Claire Rehfuss

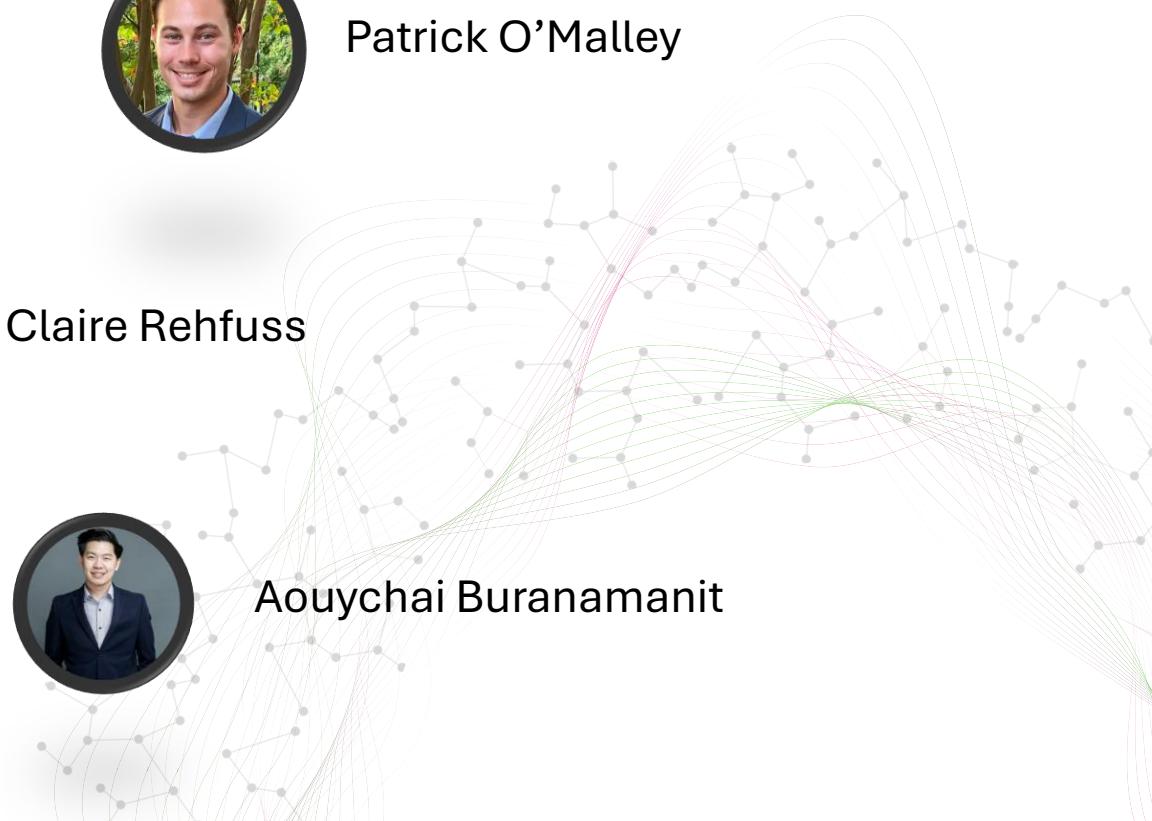


Aditya Agrawal



Aouychai Buranamanit

Thank you!

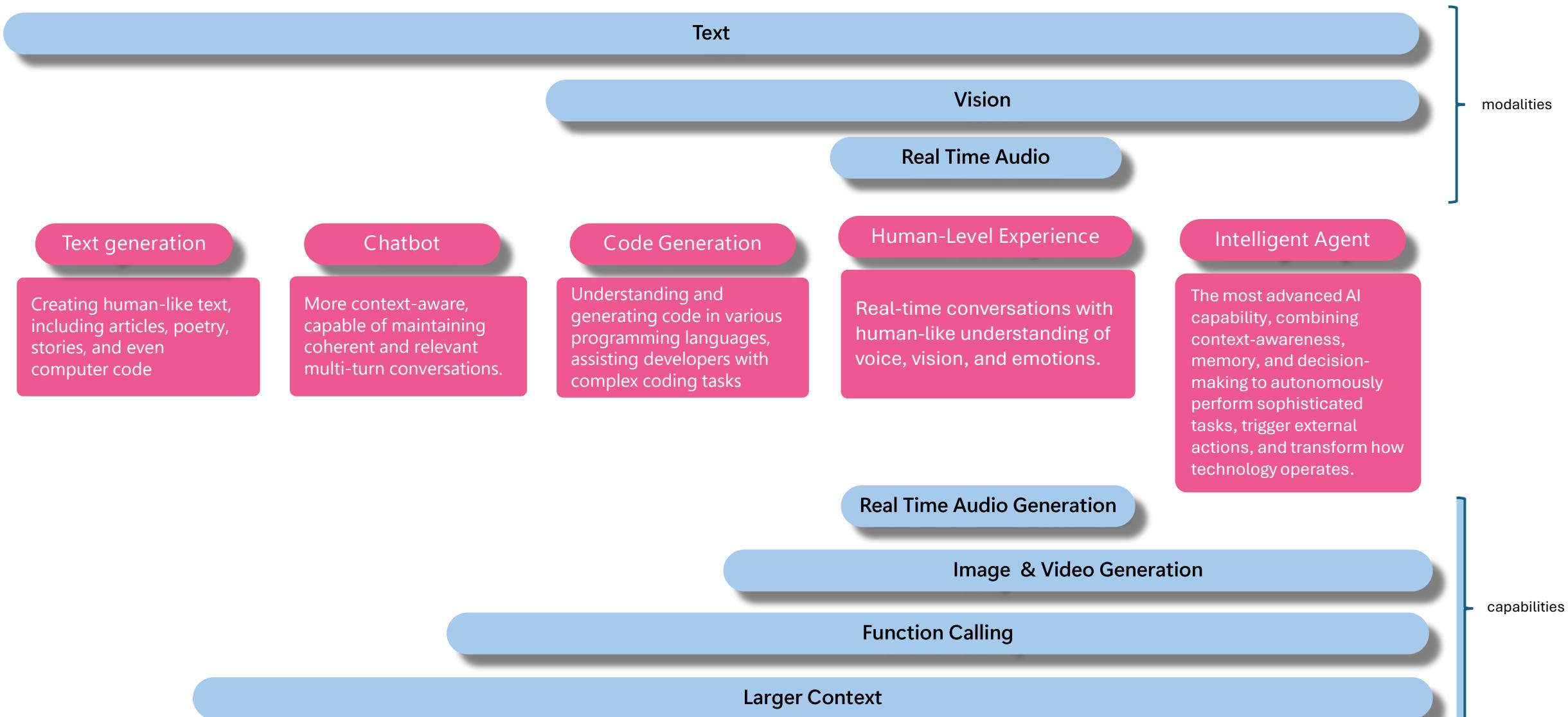




Day 1

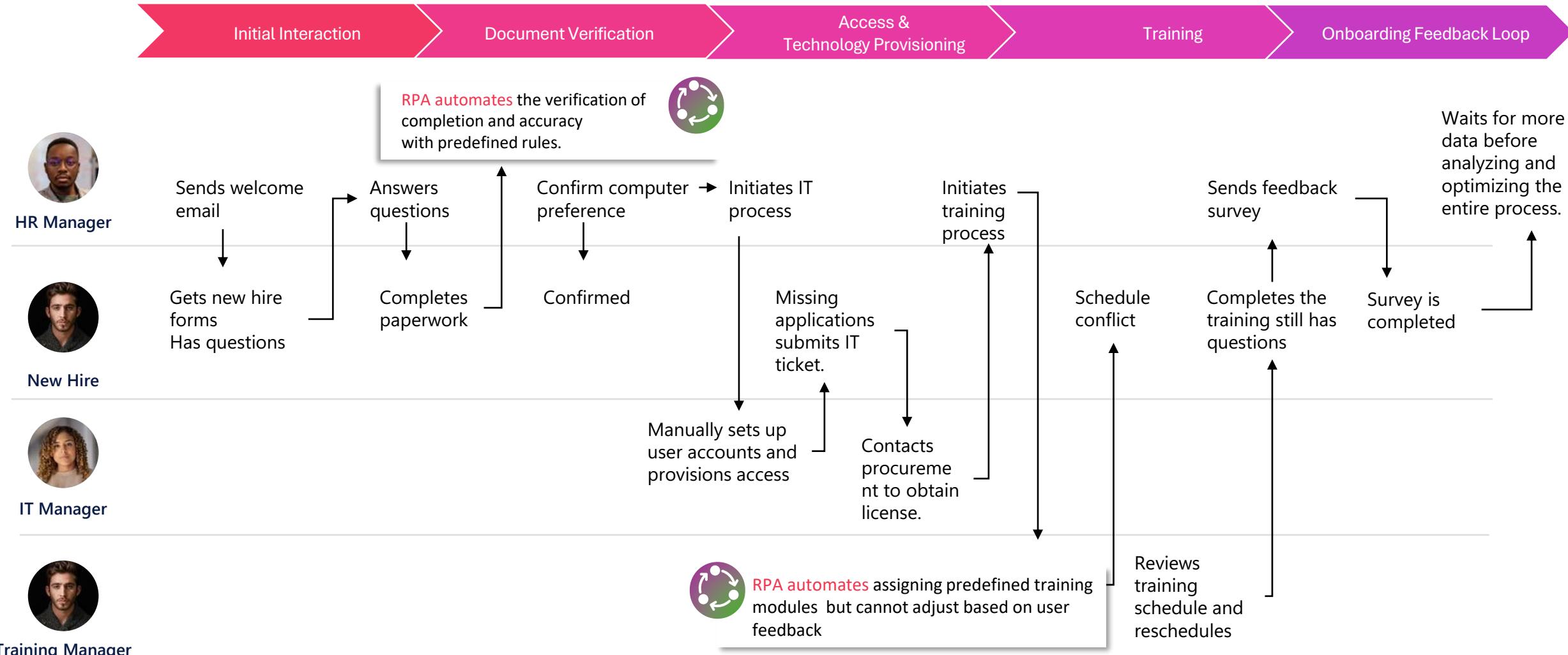
Introduction to Agentic AI

From GenAI to Intelligent Agent



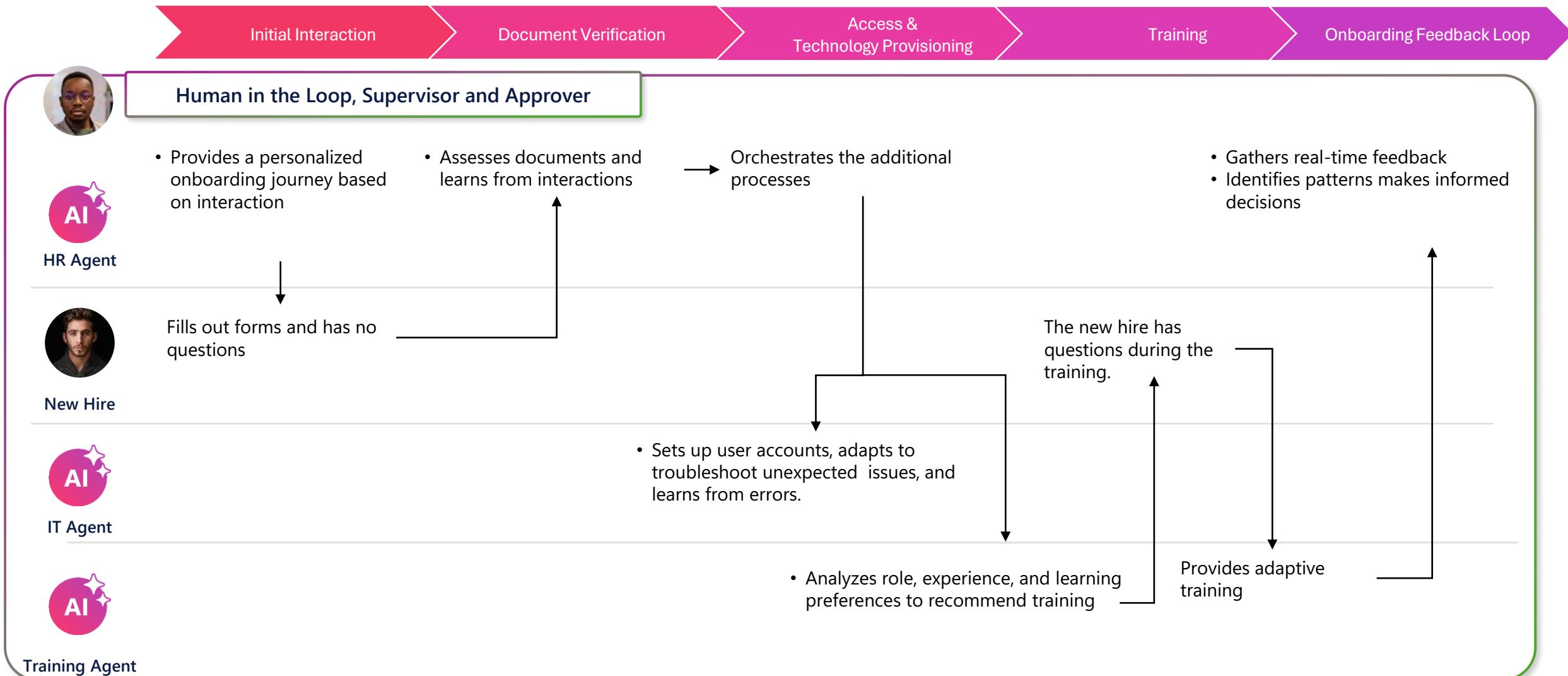
TODAY – Workflow Automation

Existing solutions can only automate very specific tasks that have clear inputs and outputs

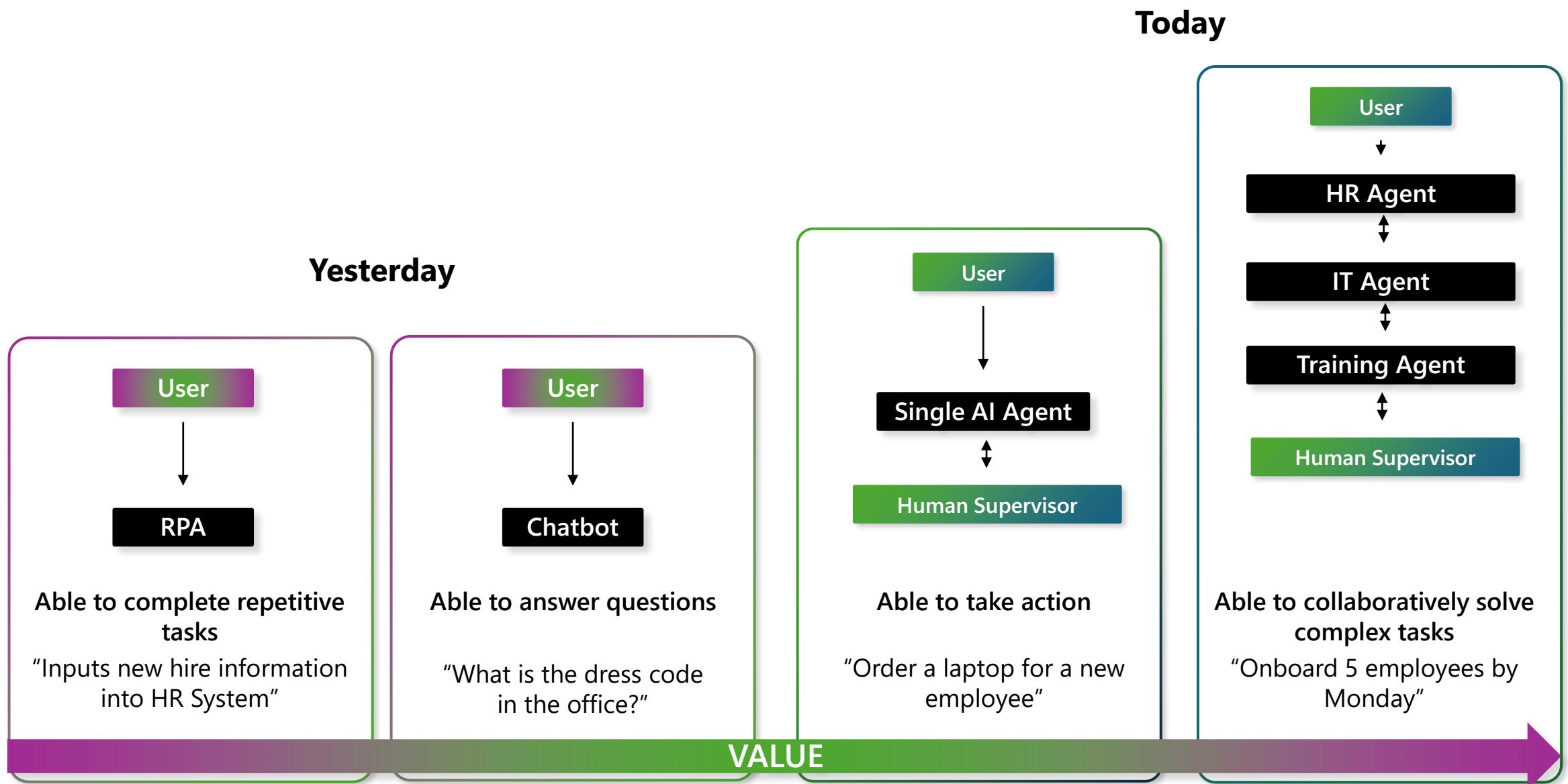


TOMORROW – AI Agents

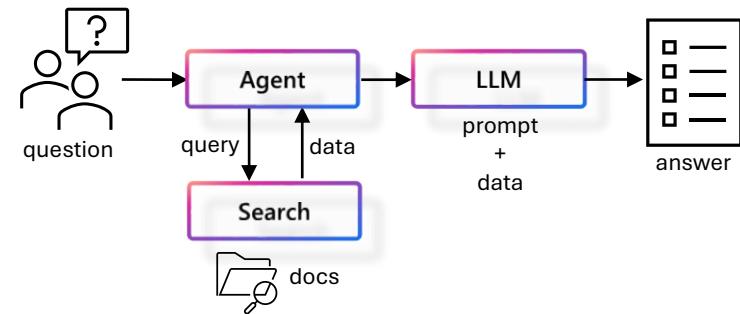
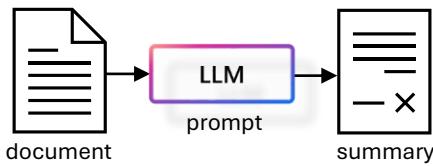
With AI Agents, these steps can be fully automated for the first time.



Promising even more efficiency, value, and advantage



Spectrum of LLM-based Solutions



No Agent

Narrow one shot task

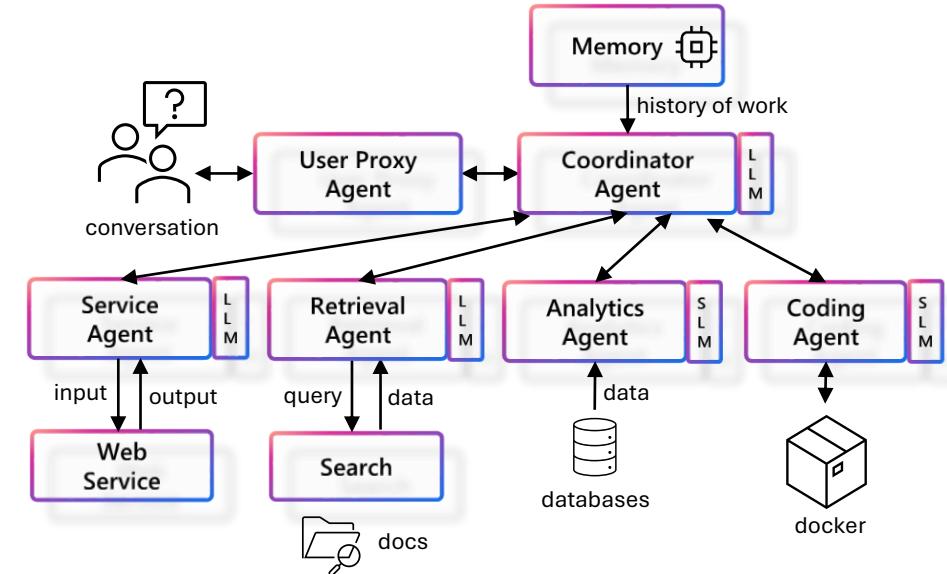
Ex: log to JSON

Single Agent

Clearly scoped iterative task

Ex: providing an answer with supporting evidence to a complex question

VALUE

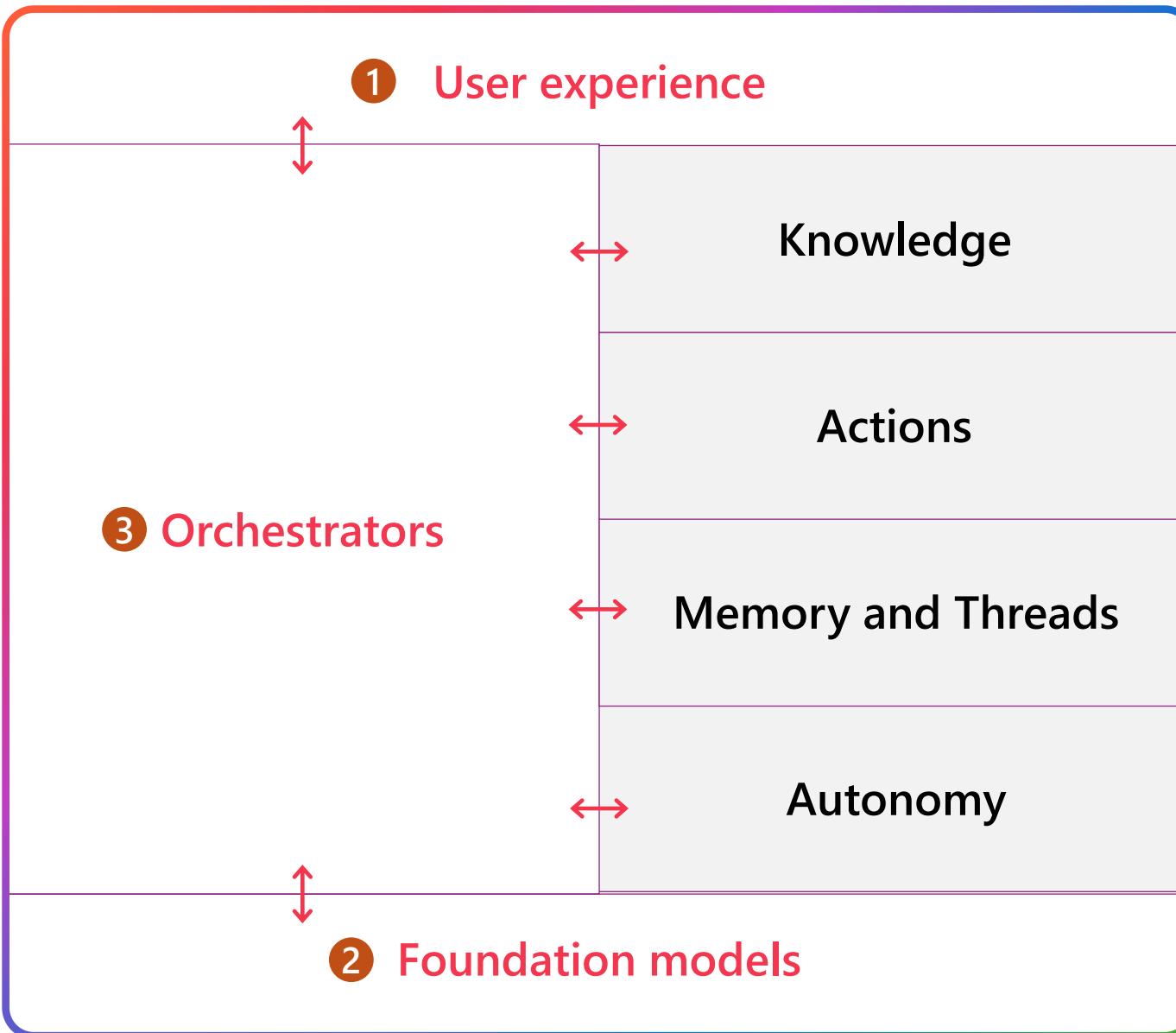


Multi-agent Systems

Wide scope complex use case requiring diverse skills

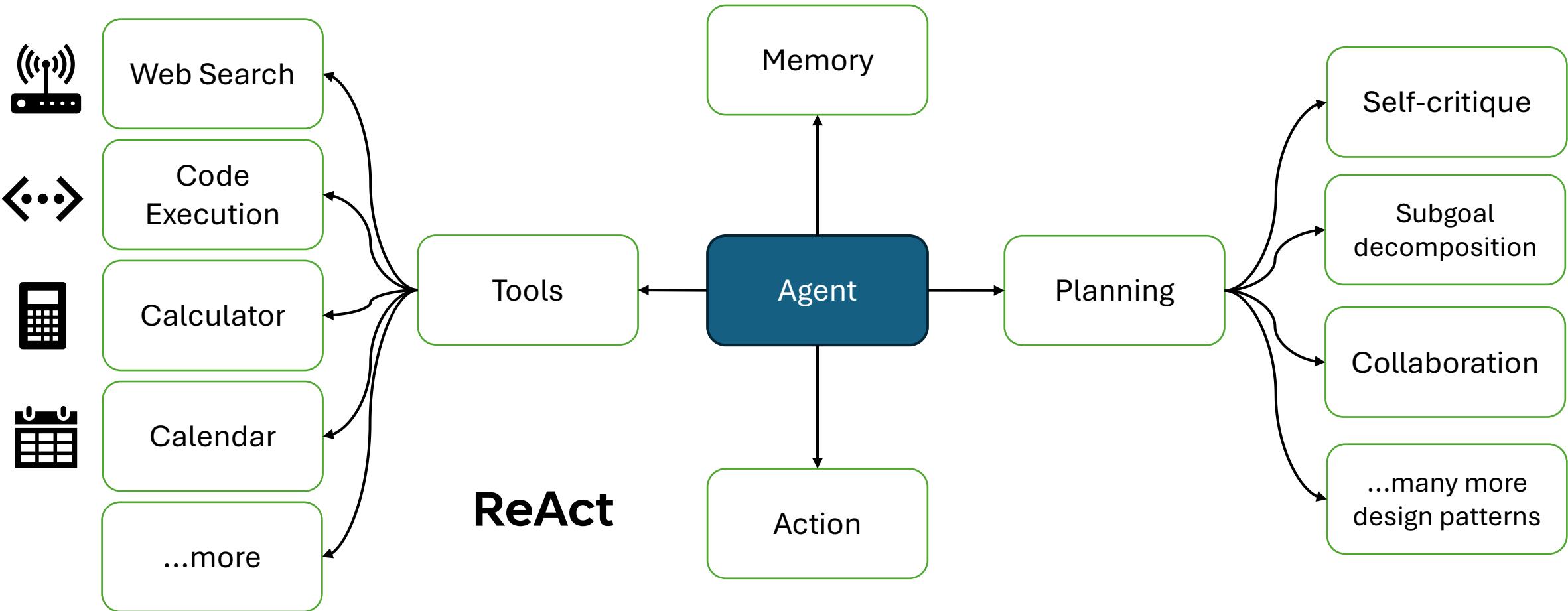
Ex: Propose 2 Instagram marketing campaigns including assets that would leverage the top 2 recent trends in our past quarter US Sales to boost our mailing list user base and predict the impact of each campaign

Components of AI Agents

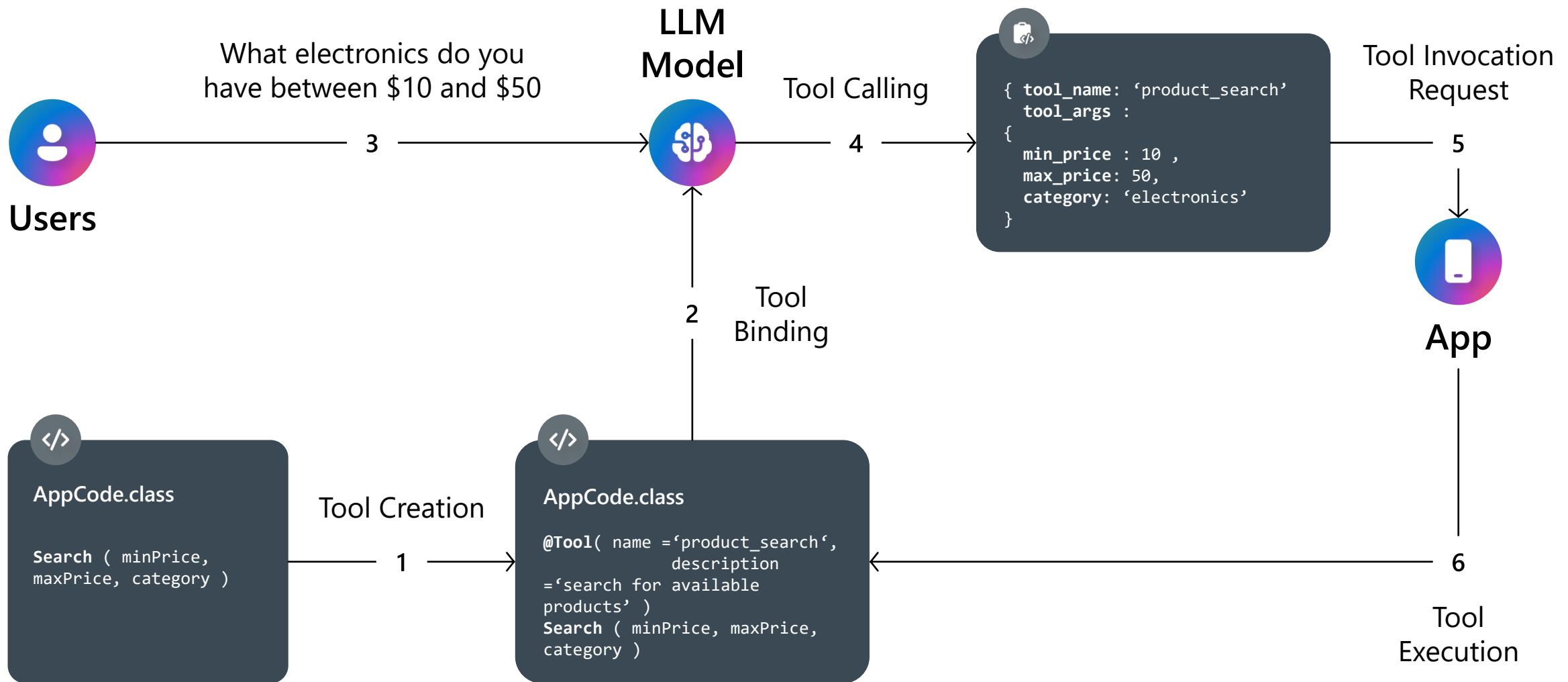


- User Experience:** The prompt that starts the whole process of the agent execution, as well as the human-agent interaction.
- Knowledge:** Search and retrieve information from online sources or company knowledge base, to ground the models.
- Actions:** Helps the agents perform certain actions (e.g. send an email, write a report) through connections to key applications.
- Memory and Threads:** Captures and stores the past interactions for hyper-personalisation and increased human-like interactions
- Autonomy:** Leaves the agents perform the tasks through an event-driven, trigger-based approach.
- Foundation Models:** Enables AI Agents to think throughout the process, helping to plan and reflect.
- Orchestrators:** Whether to be the client-side code or orchestrate across multiple agents across multiple cloud, orchestrators are key to bringing everything together.

Agentic AI capabilities



Agentic Pattern - Tools Calling



Agentic Pattern - Memory

Short Term

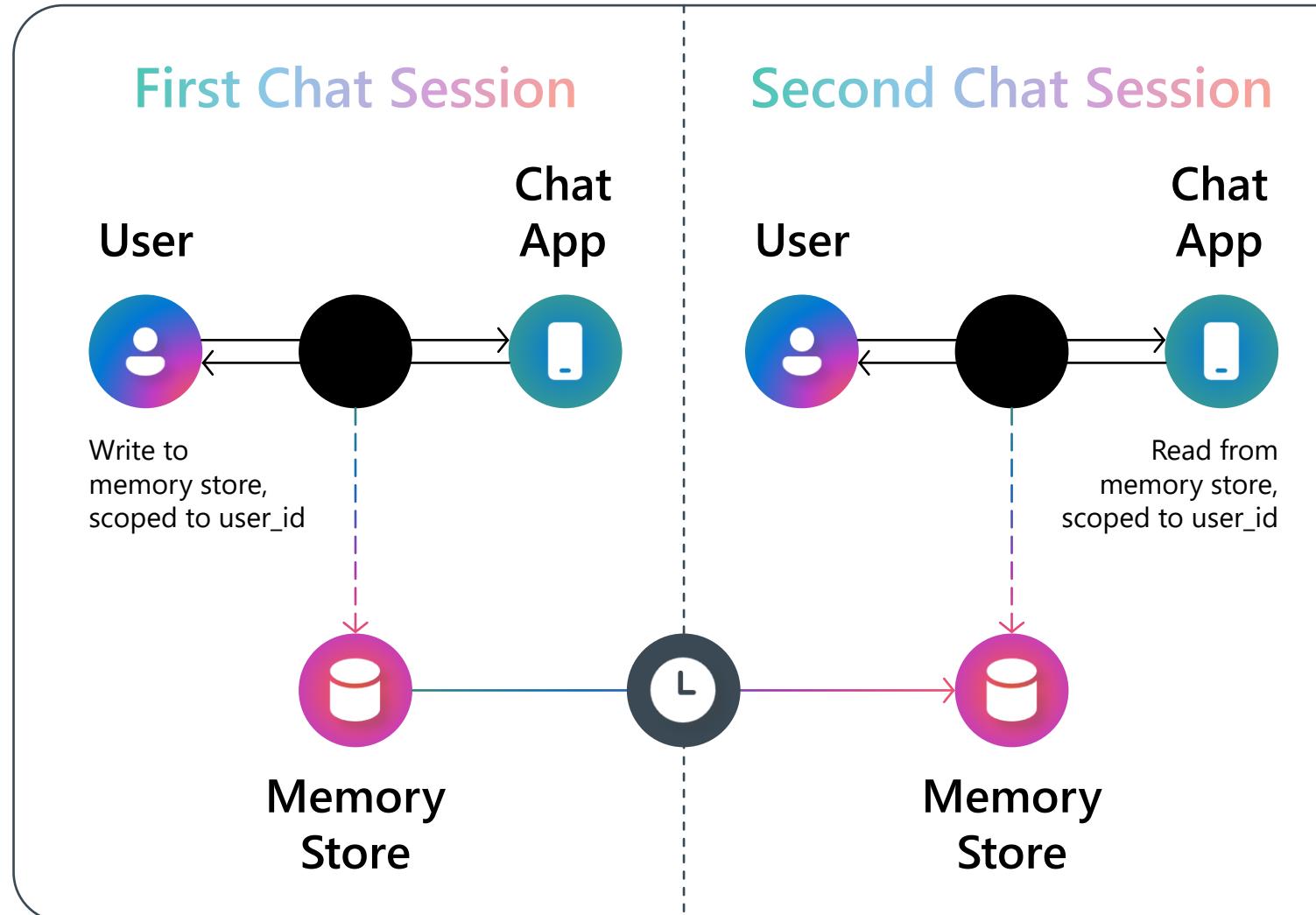
- Access steps info in one loop iteration
- Shared state context
- Chat history

Long Term

- Access steps info in long running conversation
- State persistence

Conversation History Truncation

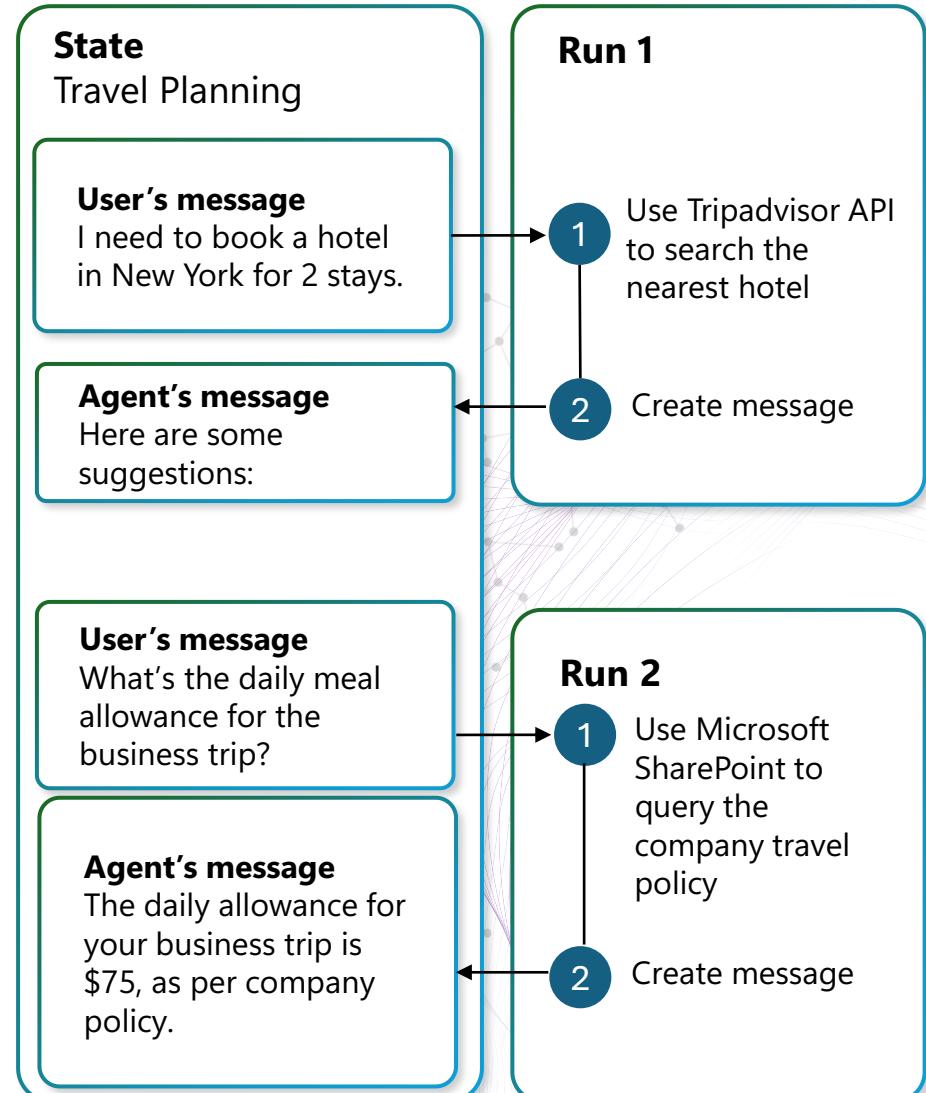
- Trim by tokens
- Trim by message count
- Trim + summary (LLM call required)



State – Providing History to Conversations

Definitions

- State/Thread – a conversation session between an agent and a user. States store Messages and automatically handle truncation to fit content into a model's context.



State Storage VS Memory

	State Storage	Memory Store
Description	CRUD ALL thread metadata, orchestration context, and conversation state in a scalable, consistent <u>operational</u> storage account	Store information about previous conversations in an index with intelligent <u>retrieval</u> capability
Requirements	Consistency, scalability for operational interactions Real-time performance	Indexing and chunking data for vector/semantic search and retrieval

Agentic System Design Patterns

Single Agent

Multi-domain

Multi-agent

Design Patterns Explored in Workshop

Intelligent Single Agent



Multi-domain Multi-agent

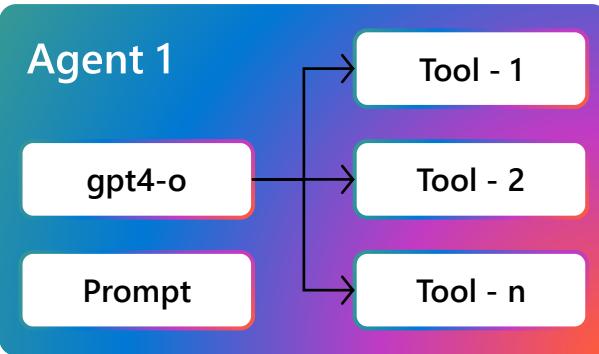


Collaborative Multi-Agent

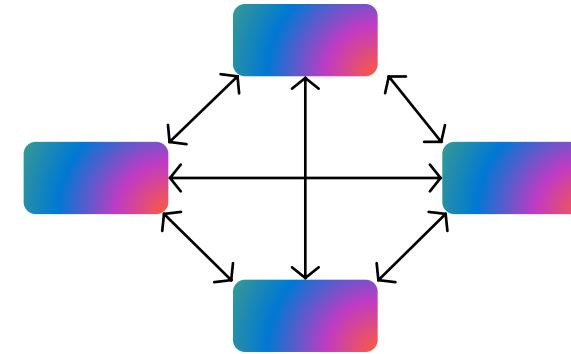


Agents orchestration and communication styles

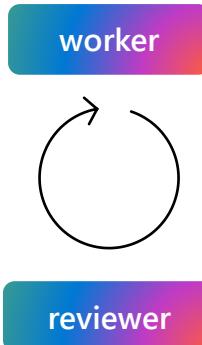
Single Agent



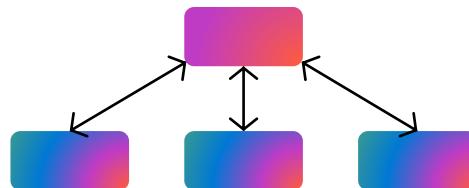
Network



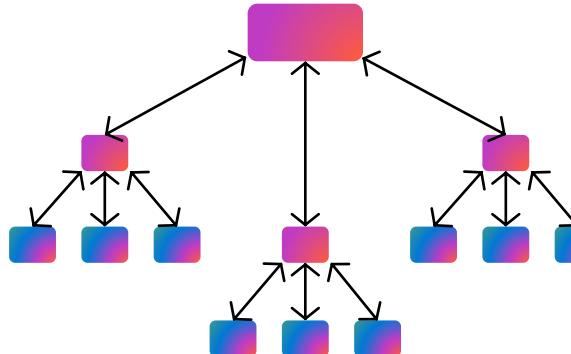
Reflection



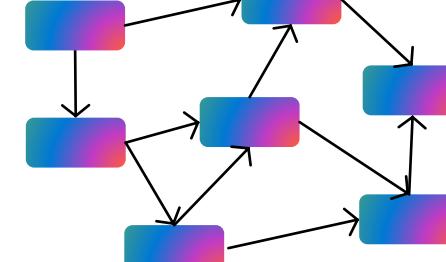
Supervisor



Hierarchical

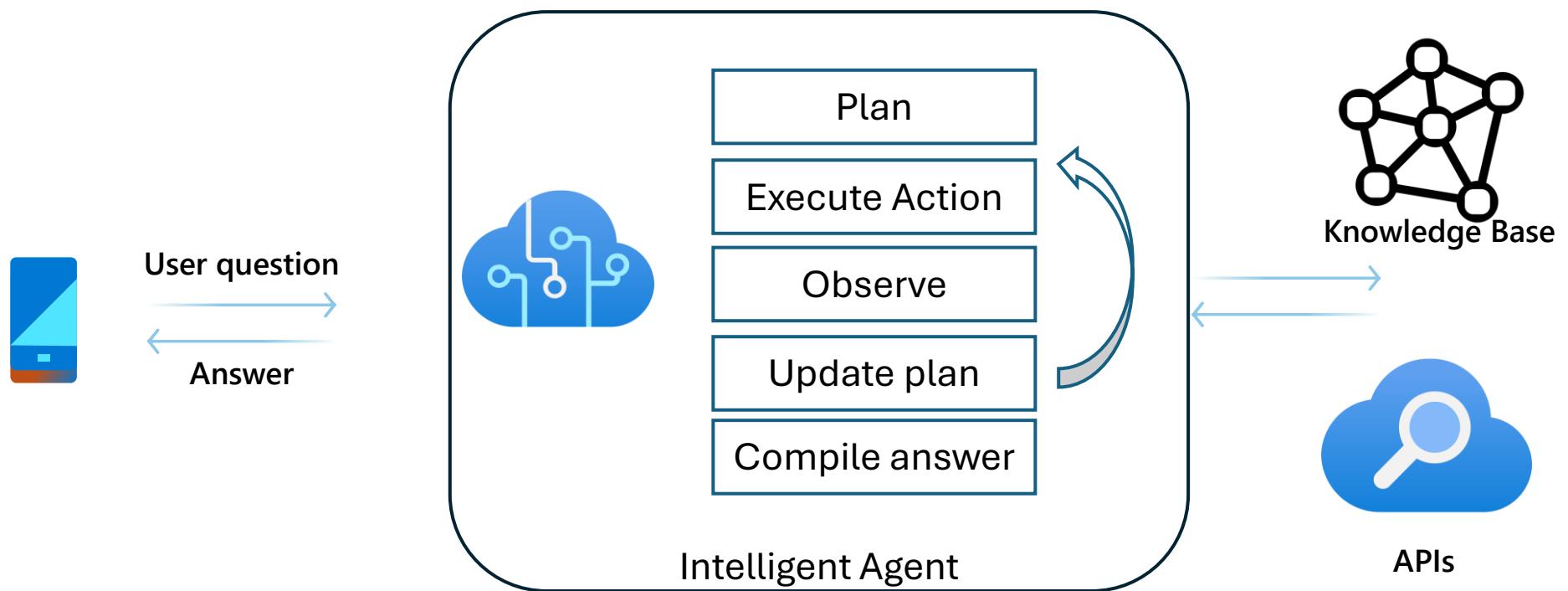


Custom



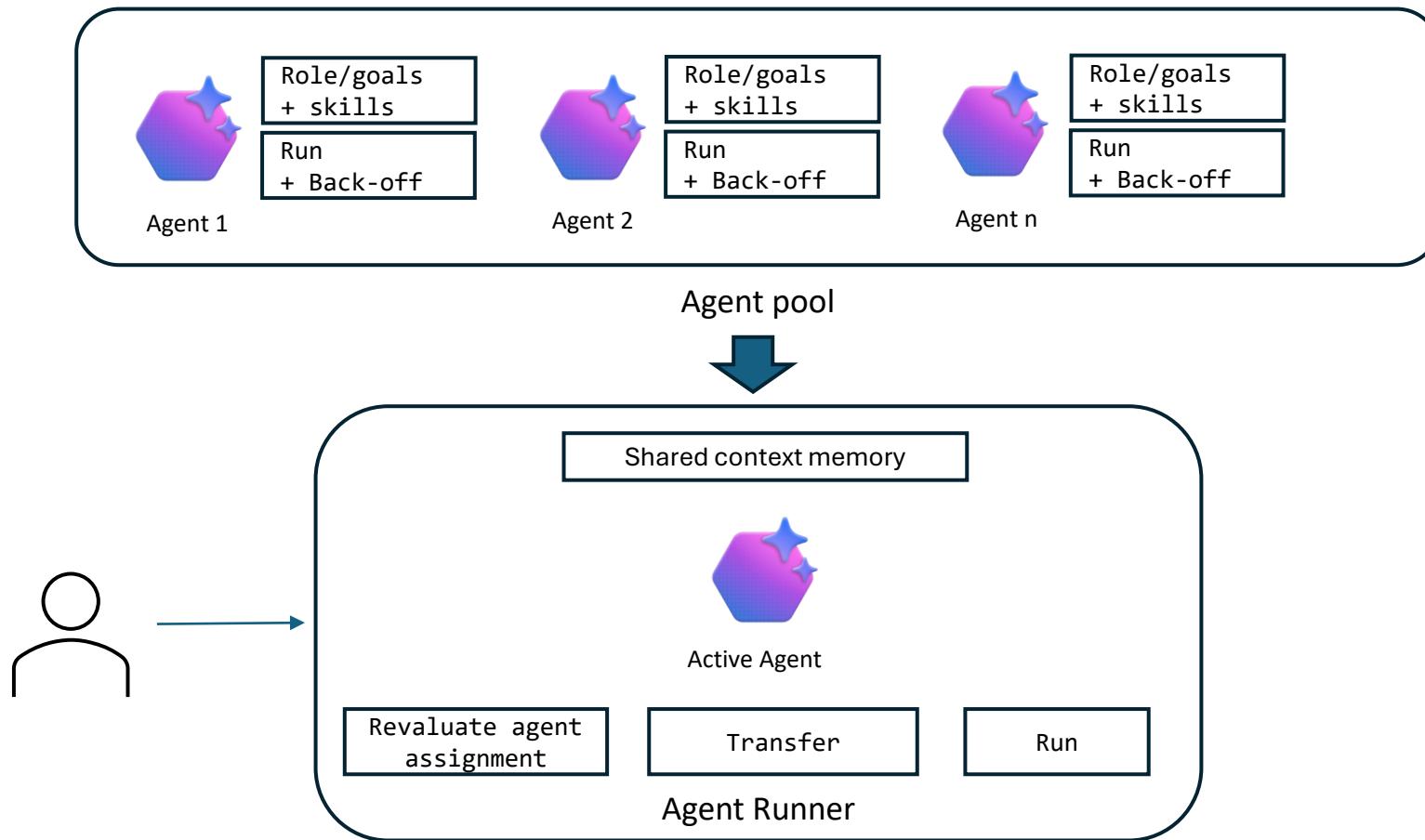
Intelligent Single Agent

Intelligent Agent with ability to translate question into a research problem and produce high quality answers



Multi-Domain Multi-Agents System

Multiple domain-specific agents are orchestrated by an Agent Runner to scale across multiple domains while appearing as a single agent to users.

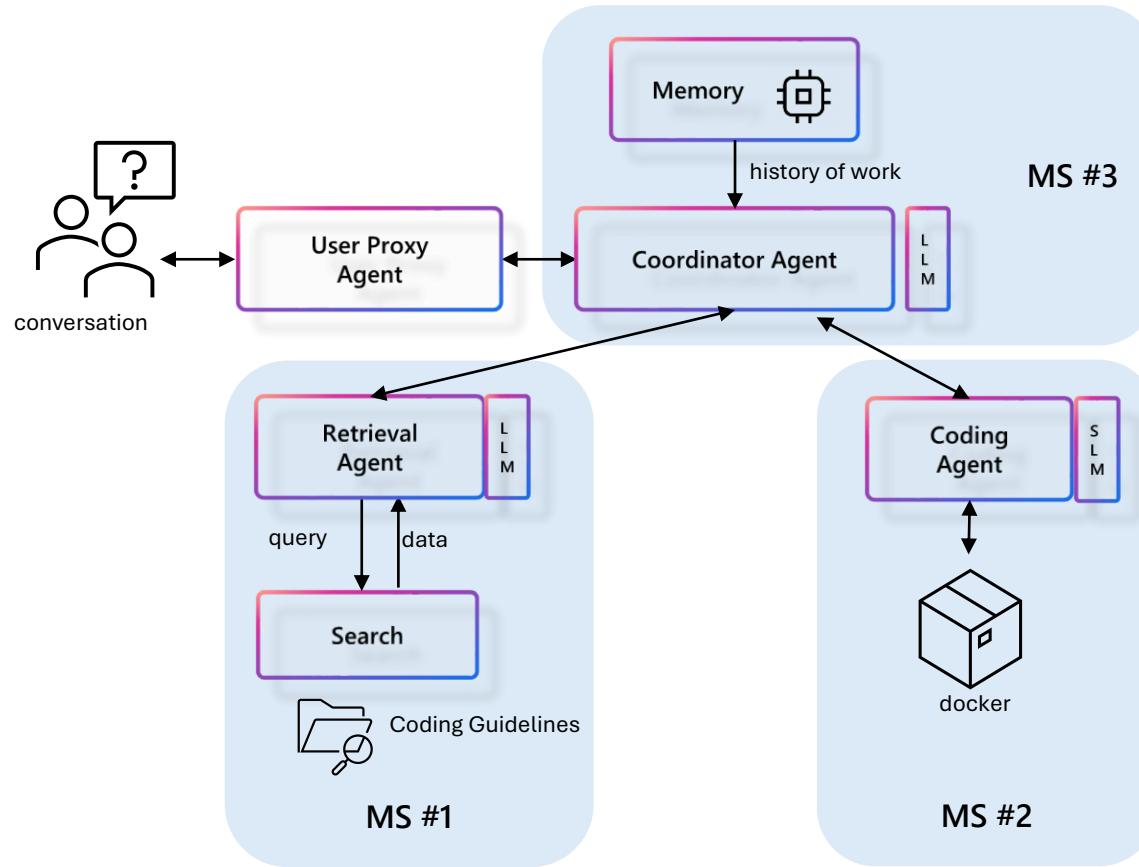


Critical Design Elements

- ✓ Agents capability descriptors
- ✓ Scalable Agent Runner able to manage 10s to 100s of agents
- ✓ Ability to manage domain switching with proper memory management
- ✓ Avoid single interceptor problem as individual agents maintain direct communication with user and can hand off when needed

Collaborative Multi-Agent System

A complex problem is decomposed into smaller, manageable parts, each addressed by specialized agents, effectively a micro-service (MS). These agents work together in a coordinated manner within a workflow to efficiently solve the overall problem.



Critical Design Elements

- ✓ Adaptive planning within scope of existing tightly scoped skills (agents)
- ✓ Handles ambiguity by discussing and refining requirements with human
- ✓ Memory to handle complex long running execution of a plan
- ✓ Effective inter agent communications
- ✓ Test, monitor, release & maintain each agent independently to quickly handle quality & safety issues

Agentic Frameworks

Azure AI Agent Service

Semantic Kernel

Autogen

Role of Frameworks

Abstraction Layers

Reduce Complexity
of integrating tools
and workflows

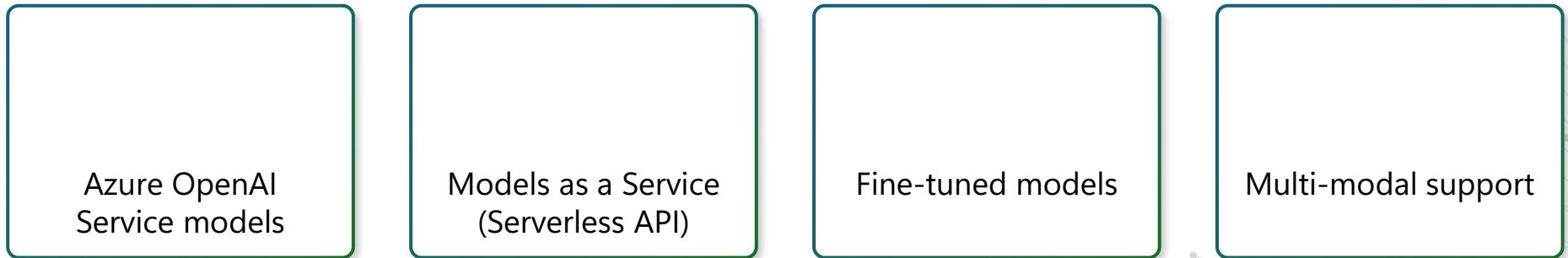
Built in Functionality

Common tools,
memory
management,
decision making

Best Practices

Streamline
development
process

Azure AI Foundry Enables Flexible Model Selection



Azure AI Studio / Model catalog

Find the right model to build your custom AI solution

Announcements

- News from Cohere!
- New SLM from Mistral
- Meta Llama 3.2 models are here!
- Experience the o1 models

All filters Collections Industry Deployment options Inference tasks Fine-tuning tasks Licenses

Models 1795

Search

Model Name	Type	Description
gpt-4o-realtime-preview	Audio generation	Real-time audio generation
openai-whisper-large-v3	Speech recognition	Large speech recognition model
openai-whisper-large	Speech recognition	Large speech recognition model
gpt-4	Chat completion	Large language model for chat completion
gpt-35-turbo	Chat completion	Large language model for chat completion
o1-preview	Chat completion	Large language model for chat completion
o1-mini	Chat completion	Small language model for chat completion
gpt-4o-mini	Chat completion	Small language model for chat completion
gpt-4o	Chat completion	Large language model for chat completion
gpt-4-32k	Chat completion	Large language model for chat completion
gpt-35-turbo-instruct	Chat completion	Large language model for instruction
gpt-35-turbo-16k	Chat completion	Large language model for chat completion
dall-e-3	Text to image	Large image generation model
dall-e-2	Text to image	Medium image generation model
whisper	Speech recognition	Medium speech recognition model
tts-hd	Text to speech	High-quality text-to-speech model
tts	Text to speech	Medium text-to-speech model
text-embedding-3-small	Embeddings	Small text embedding model
text-embedding-3-large	Embeddings	Large text embedding model
Phi-3.5-MoE-instruct	Chat completion	Large MoE language model for instruction
Phi-3-mini-4k-instruct	Chat completion	Medium MoE language model for instruction
Phi-3-medium-4k-instruct	Chat completion	Medium MoE language model for instruction
Phi-3-mini-128k-instruct	Chat completion	Large MoE language model for instruction
Phi-3-medium-128k-instruct	Chat completion	Medium MoE language model for instruction
Phi-3-small-8k-instruct	Chat completion	Small MoE language model for instruction
Phi-3-small-128k-instruct	Chat completion	Medium MoE language model for instruction
Phi-3.5-vision-Instruct	Chat completion	Large vision language model for instruction

Filter by

Collections

- Curated by Azure AI 200
- Azure OpenAI 26
- Microsoft 21
- Meta 44
- Mistral 13
- NVIDIA 5
- AI21 Labs 2
- Deli AI 4
- Nixtla 1
- JAIS 1
- Cohere 8
- Databricks 3
- Snowflake 2
- Hugging Face 1595
- SDAIA 1

Deployment options

- Managed compute 1757
- Serverless API 59

Industry

- Health and Life Sciences

Inference tasks

- Text generation 373
- Fill mask 321
- Text classification 274
- Text to text generation 239

Azure AI Agent Service



Azure AI Agent Service

Trust

Customer control over data, networking, and security

- BYO-file storage
- BYO-search index
- BYO-virtual network
- BYO-thread storage
- Tracing/monitoring
- Evaluation

Choice

Model choice and flexibility with the model catalog



Azure OpenAI Service

o3-mini, o1, GPT-4o, GPT-4o mini



Models-as-a-Service



Llama 3.1-405B-Instruct



Mistral Large, Small



Cohere-Command



DeepSeek v3

Skills

Richest set of enterprise connectivity

Knowledge



Actions



Logic Apps*



Azure functions



OpenAPI

Azure AI Foundry portal

Azure AI Foundry SDK

AIAS Enterprise Readiness



Bring your own storage



Keyless setup
and
authentication



Private Virtual
Network support



Tracing/
monitoring



Content filters

Try the new Azure AI Agent Service experience

	Existing	Azure AI Agent Service Experience
Private virtual networks and limitless scaling and agent monitoring	✗	✓
Securely ground agents in Bing, SharePoint, Fabric, and Azure AI Search knowledge sources.	✗	✓
Automate complex workflows through powerful action connectors	✗	✓
Leverage pre-built tools for data analysis and retrieval augmented generation (RAG)	✗	✓
Access models Microsoft Research, Hugging Face, Llama, Mistral, DeciAI, and more via hosting or inferencing	⚠ View only	✓
Scalable, secure memory management	✗	✓

[Continue with existing](#) [Create a new project](#)



Semantic Kernel

Semantic Kernel is lightweight, **open-source**, production-ready, orchestration middleware that lets you easily add AI to your apps.

- Full SDK designed to build AI agents with ease, excellent for single agents and can be extended for multi-agents with integrations to AutoGen
- Extensible and compatible with models LLMs or SLMs.
- Ideal for developers looking to leverage AI orchestration patterns similar to those used in Microsoft's Copilot systems in their own applications

1

Single-agent

Deploy agents with
Azure AI Foundry



Managed agent
micro-services

2

Multi-agent

Orchestrate them together with
AutoGen and **Semantic Kernel**



State-of-the-art
research SDK



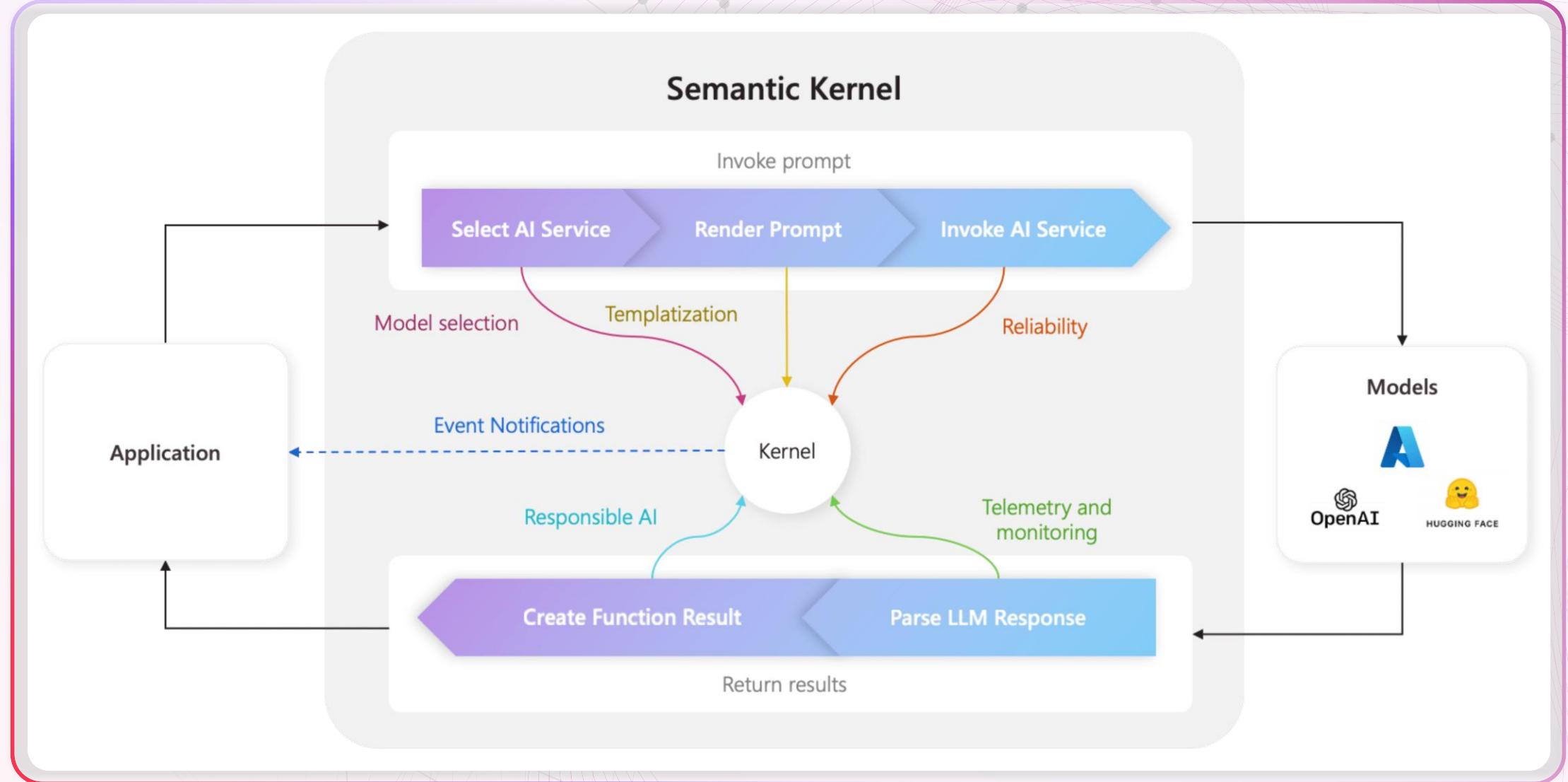
Production-ready
and stable SDK

Ideation

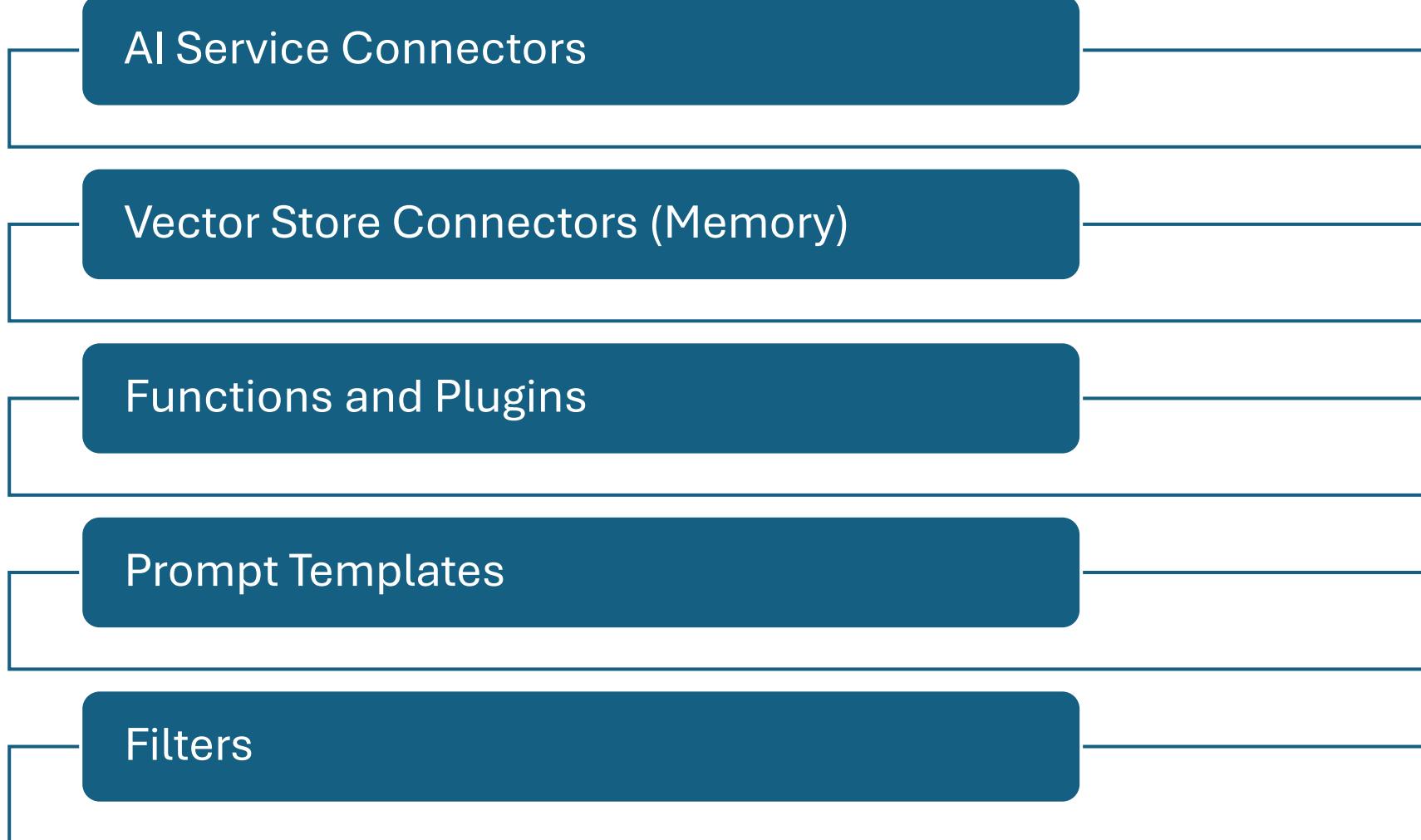
Production



The kernel is at the center



Semantic Kernel Components





What is a Plugin?

Chatbots are *nice*, but they aren't *useful* to your users until they can interact with the real world by...

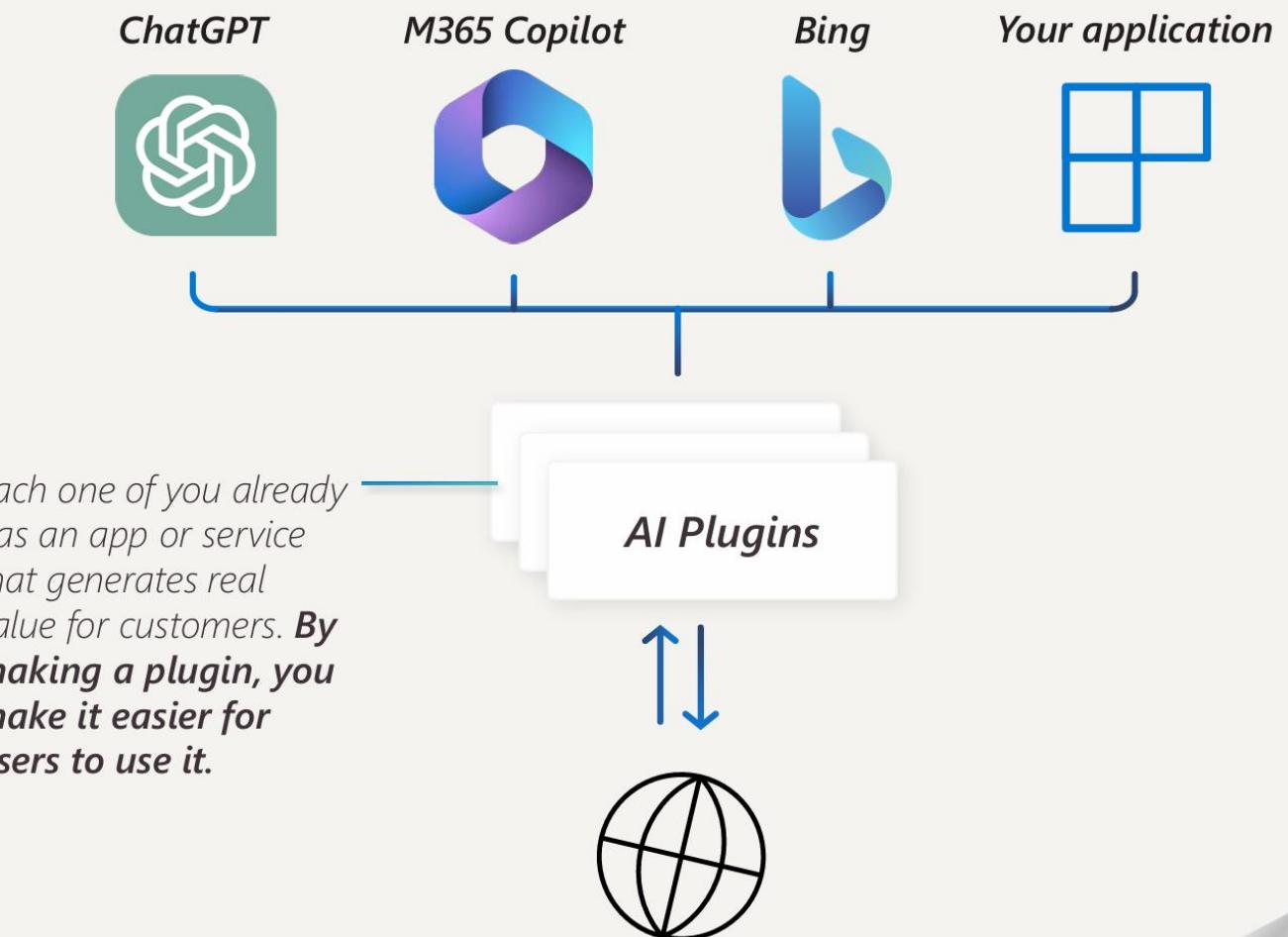
Retrieving data

Sending emails

Completing sales

Making orders

And more!



Writer plugin

Function	Description for model
Brainstorm	Given a goal or topic description generate a list of ideas.
EmailGen	Write an email from the given bullet points.
ShortPoem	Turn a scenario into a short and entertaining poem.
StoryGen	Generate a list of synopsis for a novel or novella with sub-chapters.
Translate	Translate the input into a language of your choice.

Can you write me a short poem about living in Dublin, Ireland and then create a story based on the poem?



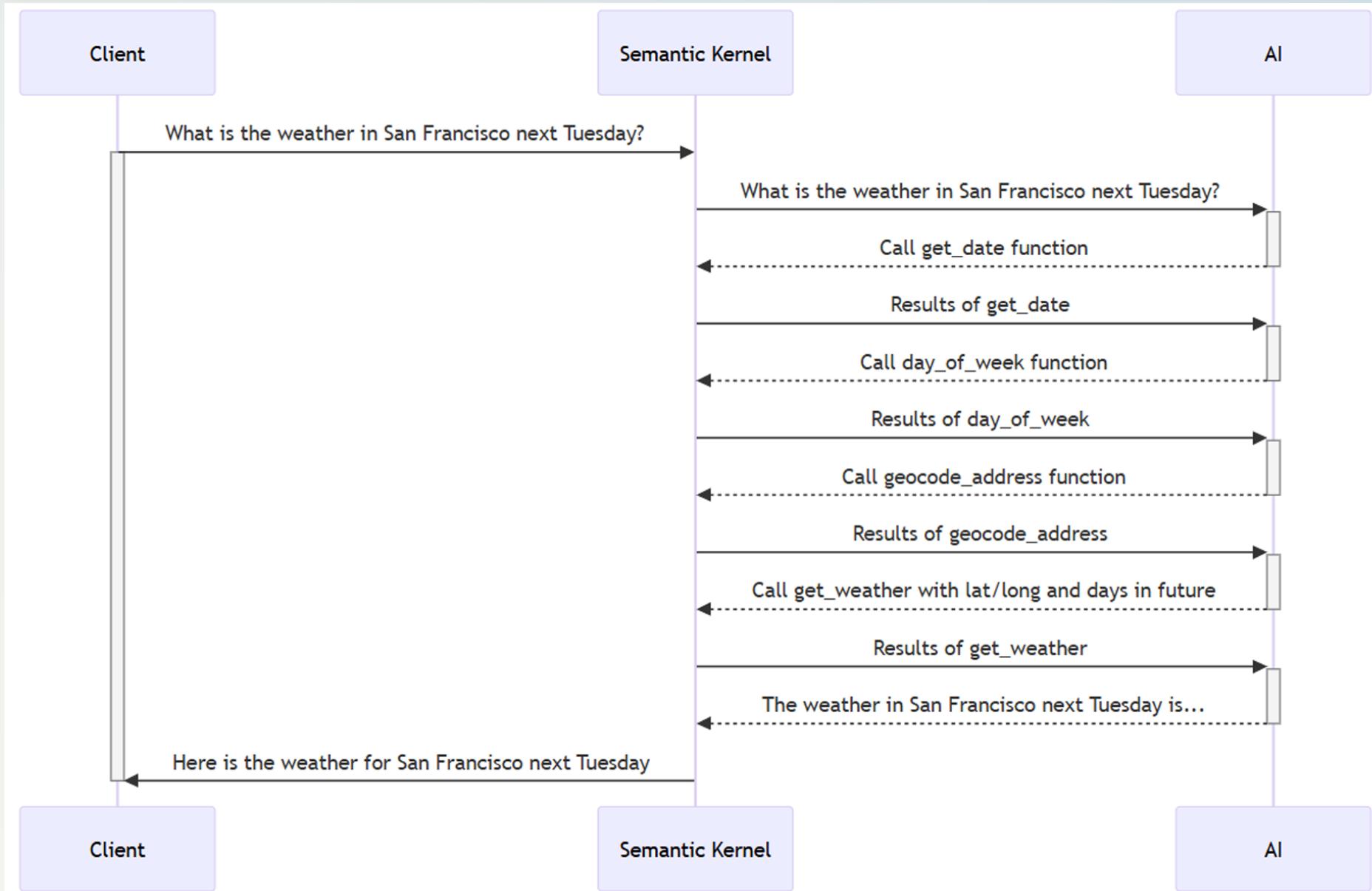
Planner

Copilot

Sure! Here's a story based on living along the Grand Canal in Dublin, Ireland...



Plugin Workflow – Follow the Magic



Semantic Kernel Agent Patterns

Single Agent

Multi-domain

Multi-agent



Autogen

- Powerful multi-agent framework with prebuilt conversation orchestration patterns for handling complex agent systems
- Extensible and compatible with models LLMs or SLMs.
- Powered by collaborative research studies from Microsoft, Penn State University, and University of Washington.
- AutoGen simplifies the orchestration, automation, and optimization of a complex LLM workflow.

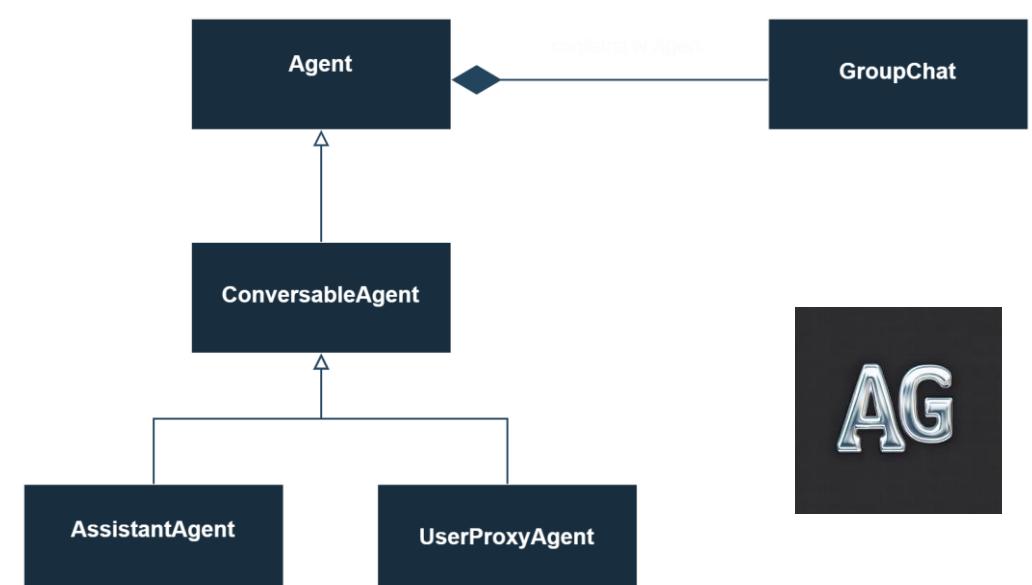
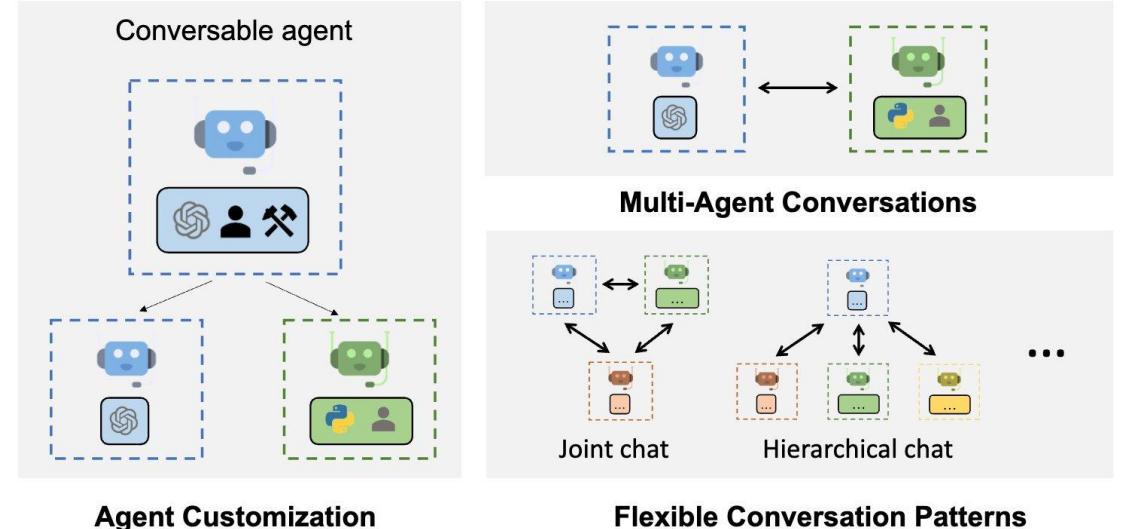
AutoGen concepts

- **Customizable and conversable agents:**

AutoGen uses a generic design of agents that can leverage LLMs, human inputs, tools, or a combination of them

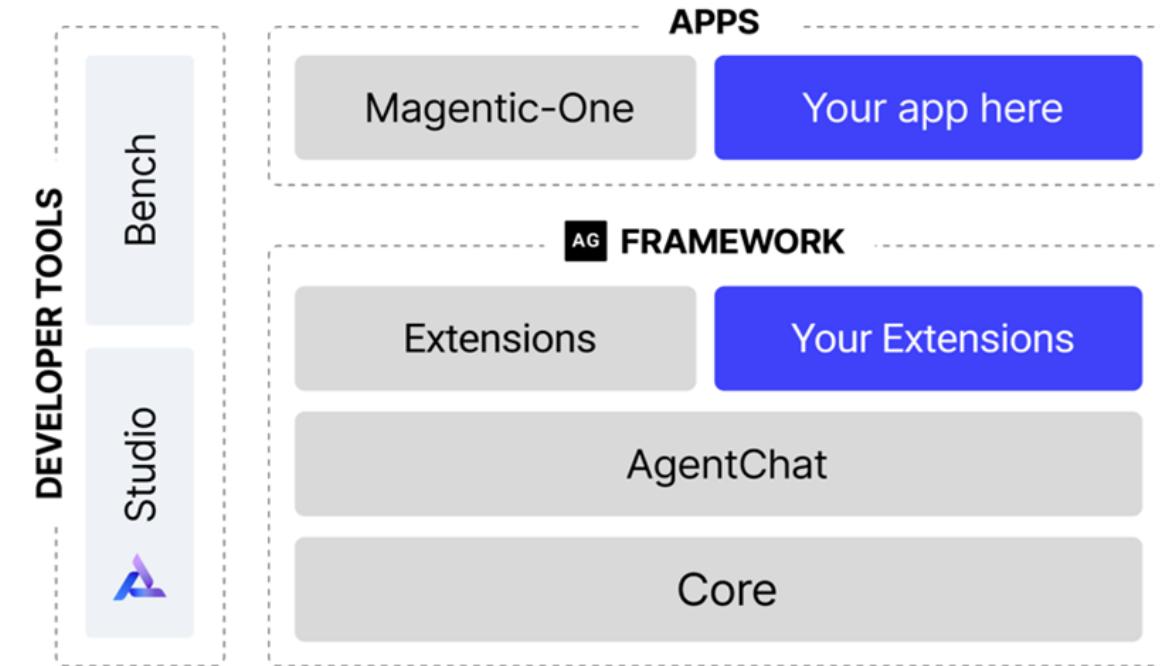
- **Conversation programming:**

- Defining a set of conversable agents with specific capabilities and roles
- Programming the interaction behavior between agents via conversation centric computation and control.

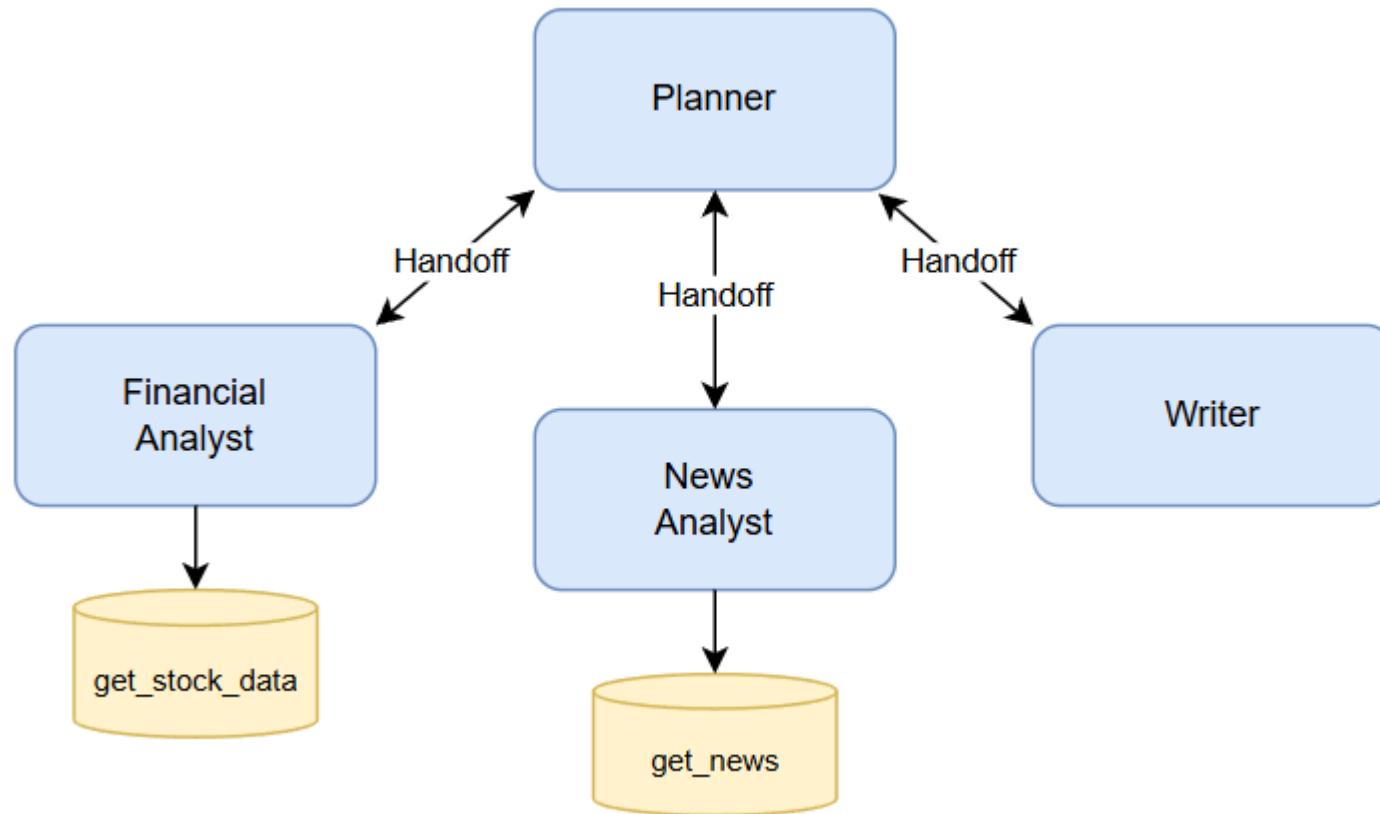


AutoGen concepts

- **AgentChat:**
High level API , with predefined roles
- **Core:**
Base framework, event-driven
programming model
- **Extensions:**
Built-in component implementations
for models, agents, tools, etc.



Stock Research Example

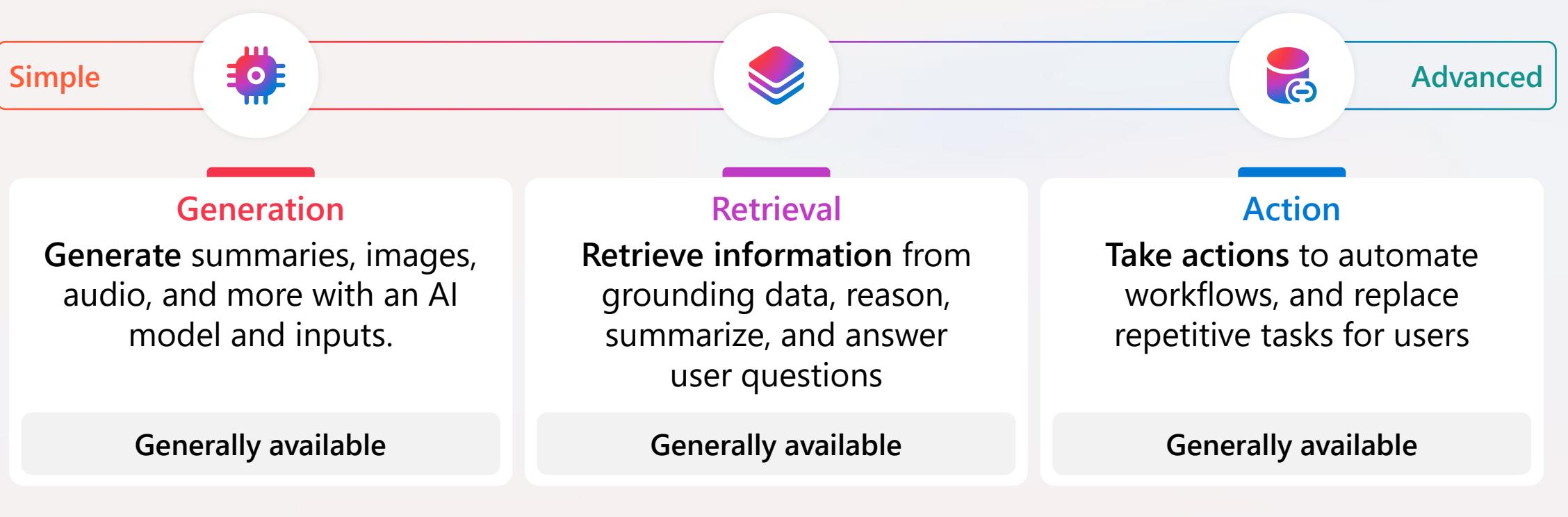


Backup slides

What are agents?

AI designed to perform a task

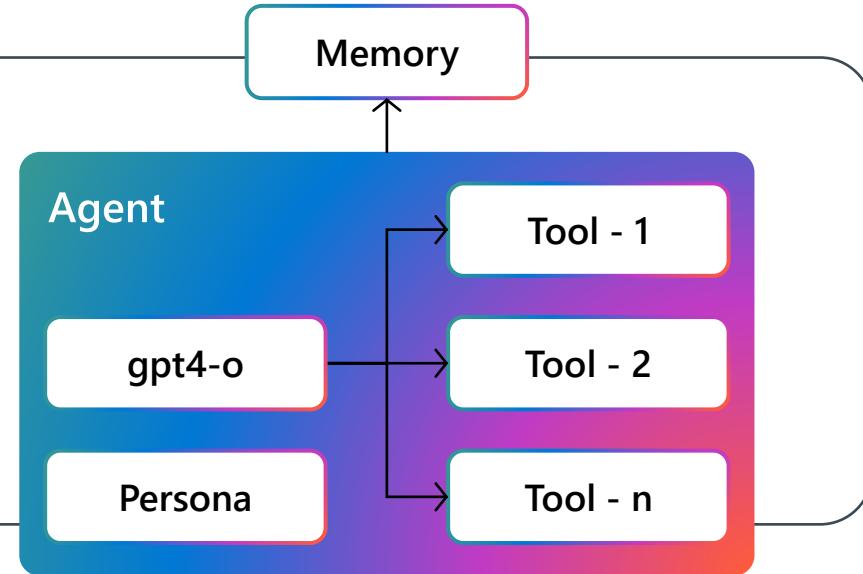
Tasks can vary in level of complexity and capabilities depending on your need



Agent Abstractions - Agent First-Class Citizen

Agent as high-level abstraction

- LLM (gpt4-o, o1 etc.)
- Persona (system prompt)
- Tools (function code calls)



Agent Chat as layer for collaboration

- Multiple agents can engage with each other
- Enables multi-turn or parallel execution



Memory - Providing Context for Agent

Memory is a **foundational** capability that allows agents to store and leverage past interactions to deliver personalized, context-aware experiences and enhance workflow efficiency.

Definition

- Memory store – the place of storage for information across multiple threads
- Memorizing – act of taking things from a thread and storing it elsewhere
- Recalling – the act of retrieving information from memory

Agent

LLM + INSTRUCTION

MEMORY

TOOLS

KNOWLEDGE

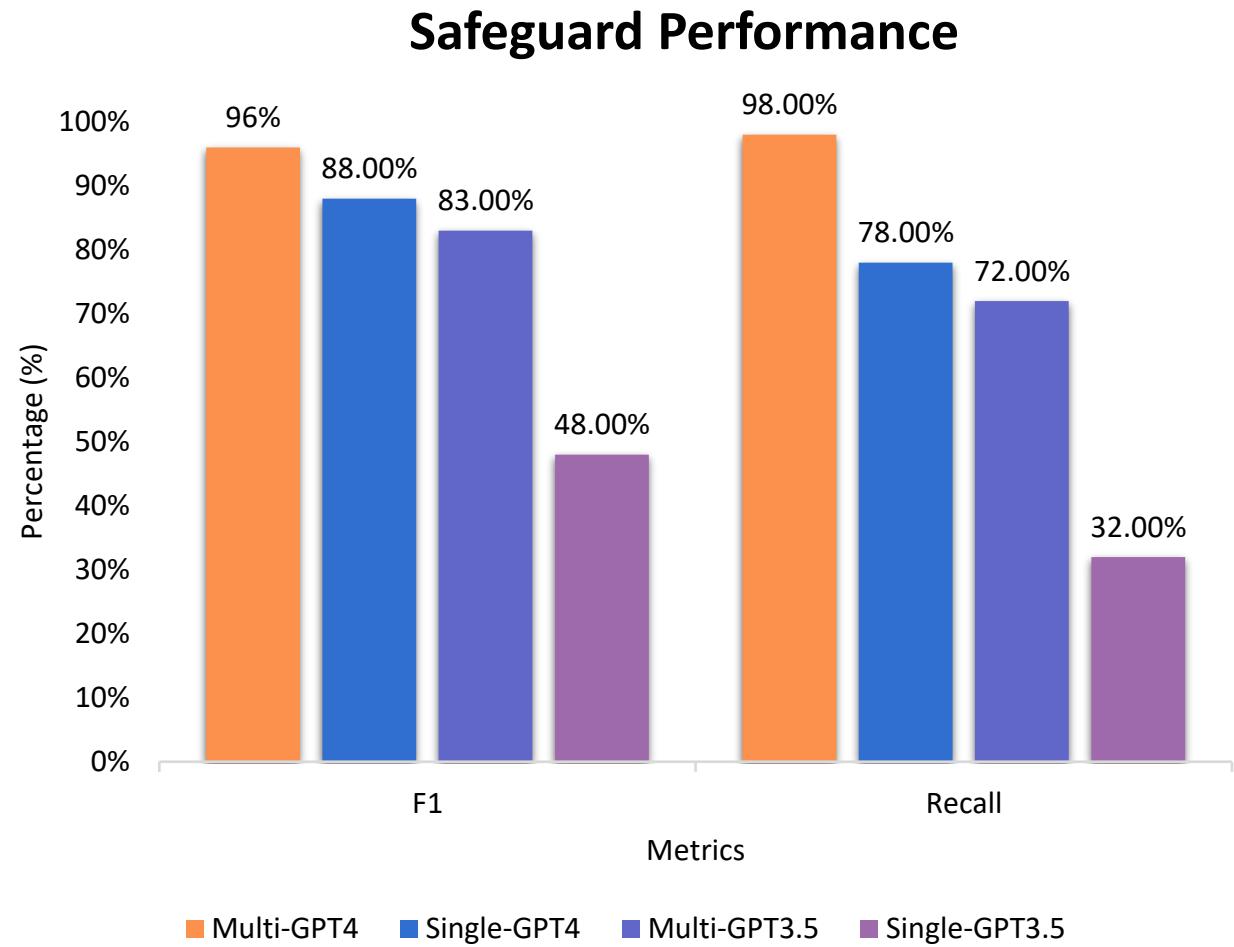


When should I use multi-agents over a single agent?

Multi agent systems can solve more complex tasks

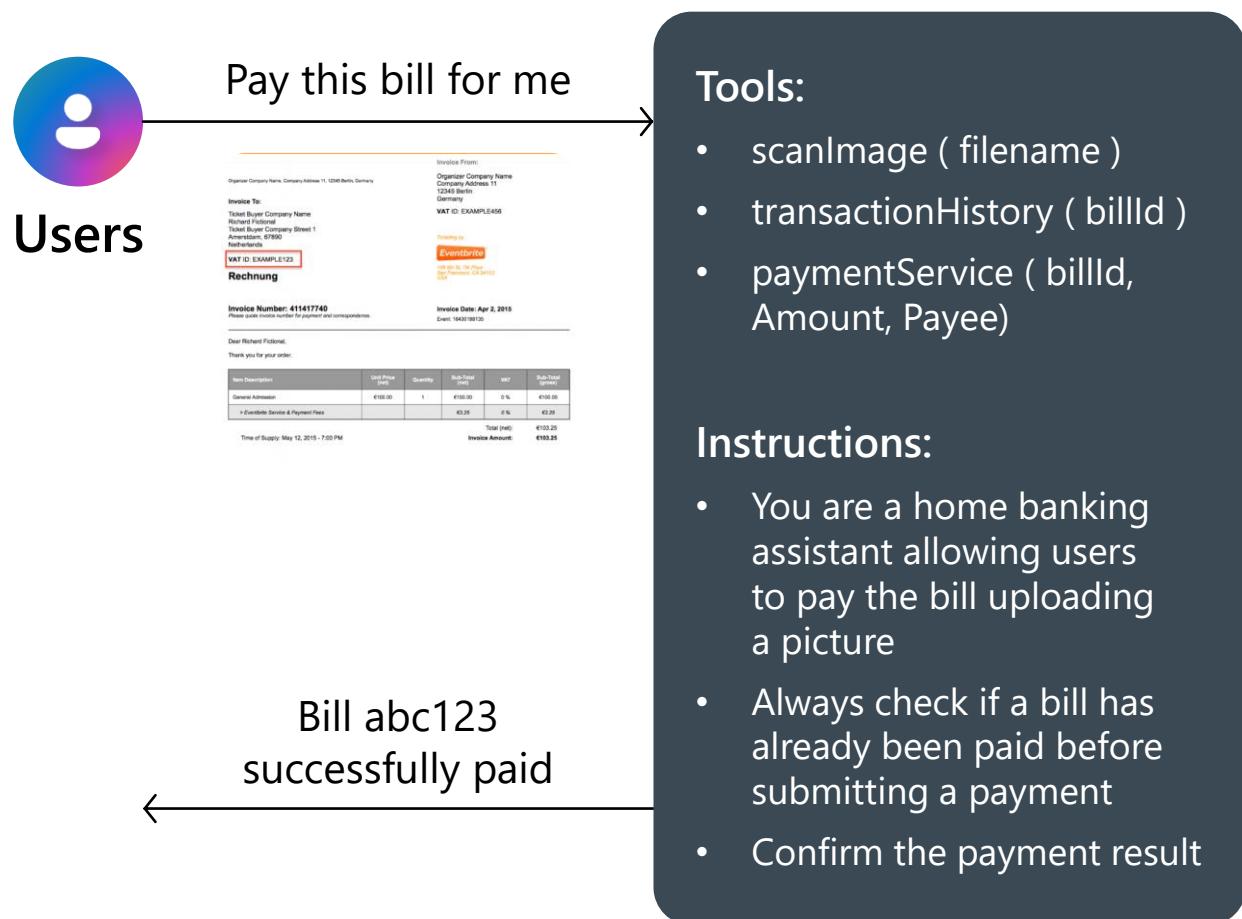
Single agents can solve a wide range of problems and are simpler to implement.

- Multi-agent systems should be used
- if a single agent is unable to solve the challenge
 - to handle tasks that involve more data, diverse roles, or complex workflows

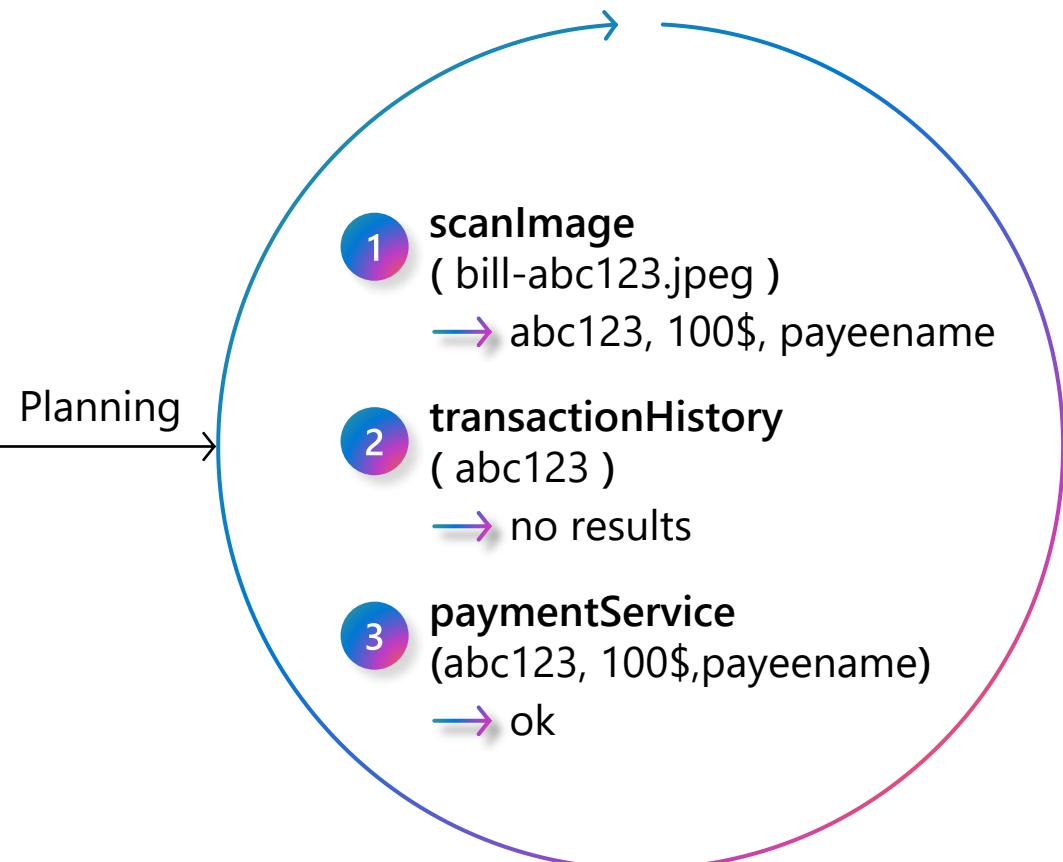


Agentic Pattern - ReAct Planning with Tools

Calling Payment Agent



`while (new tools execution request)`



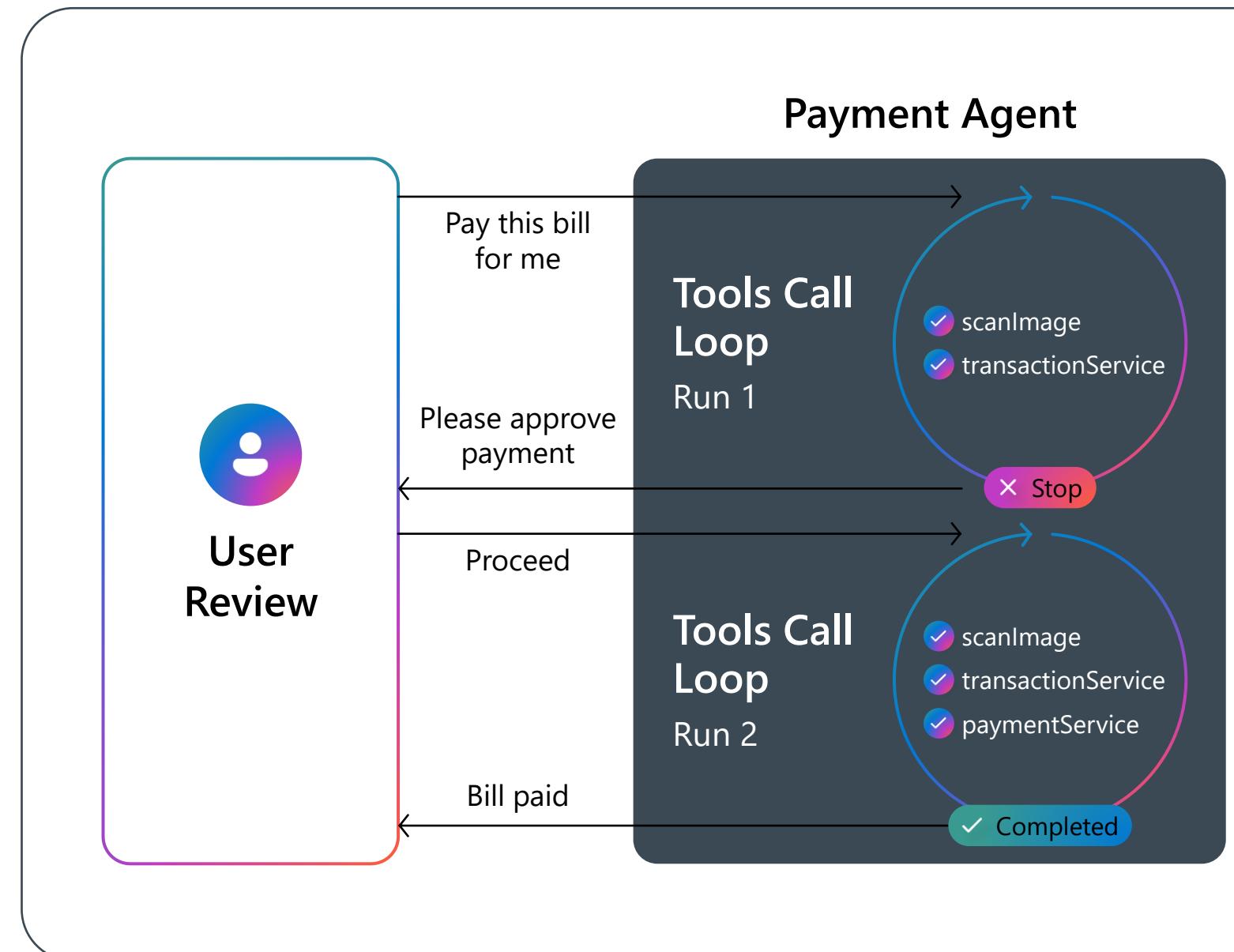
Agentic Pattern - Flow control

Looping Termination

- MaxIterations
- Message termination
- Human step /Human in the loop

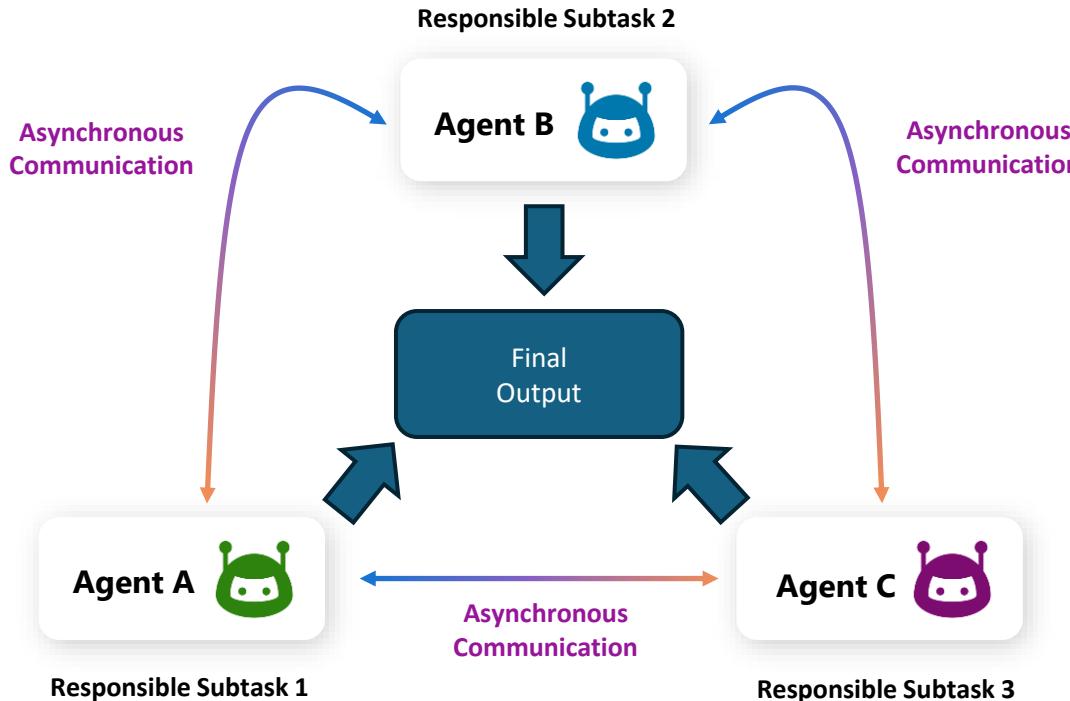
Human in the loop

- Action execution approval
- Escalation
- Data review



Cooperative vs Competitive Agents

Cooperative based-learning



Cooperative agents work together towards shared goals, enhancing problem-solving capabilities and efficiency through collaborative efforts, characterized by trust and open communication.

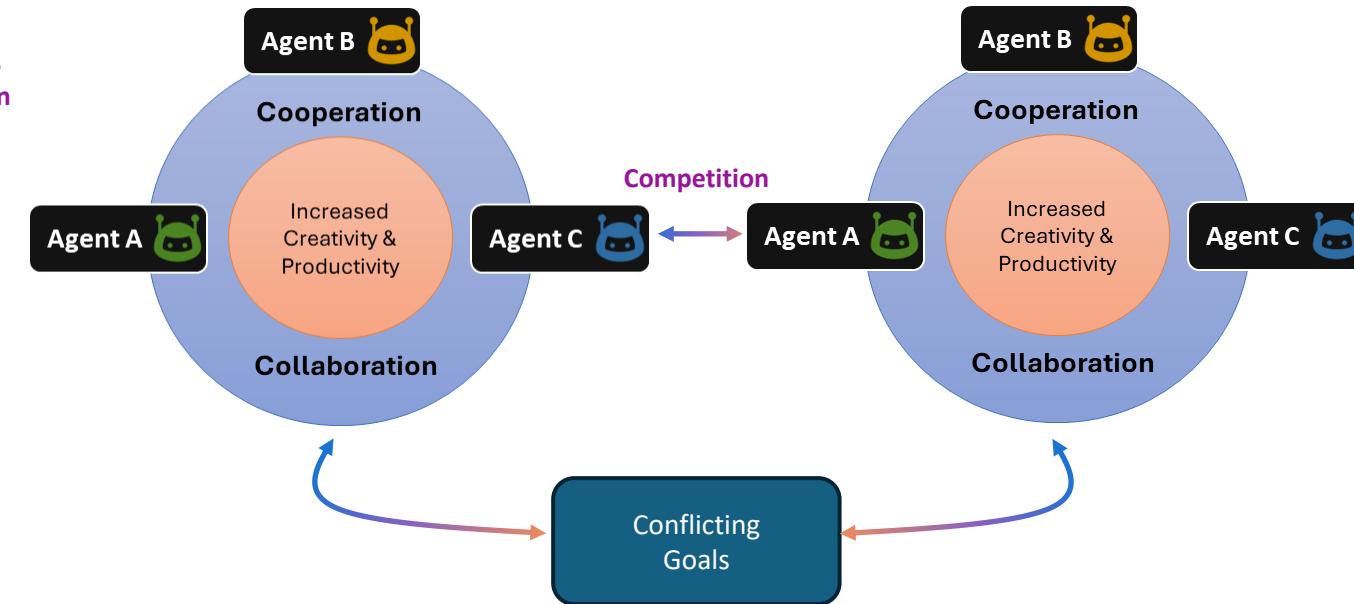
Complex Problem Solving

Efficiency

Adaptability

Competition based-learning

Seems not useful in a machine scenario



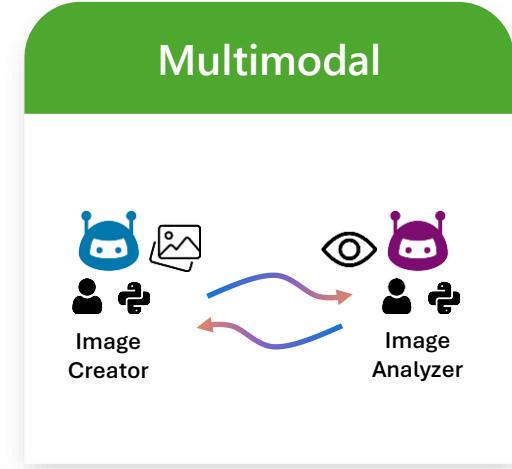
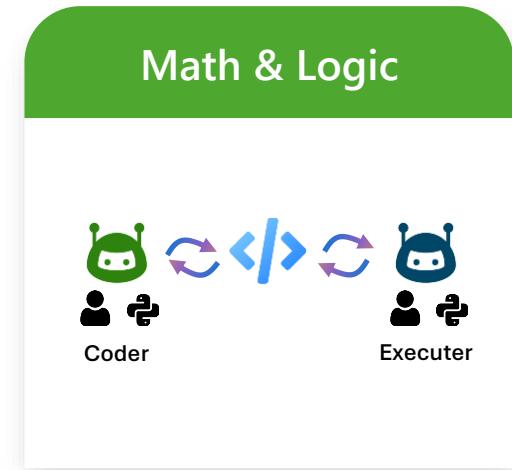
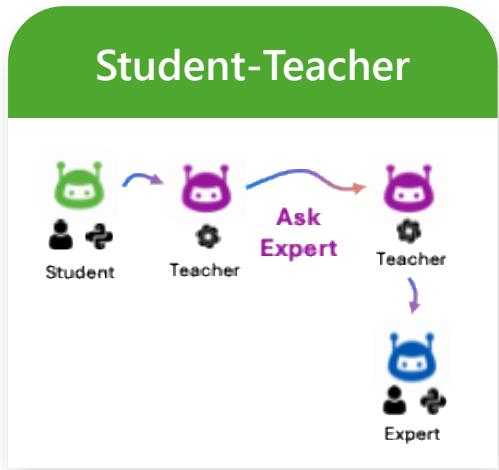
In contrast, competitive agents operate with individual goals in adversarial settings, focusing on optimizing their own outcomes, often with strategic communication and planning.

Resource Allocation

Decision Making

Realism

Multi Agent Examples



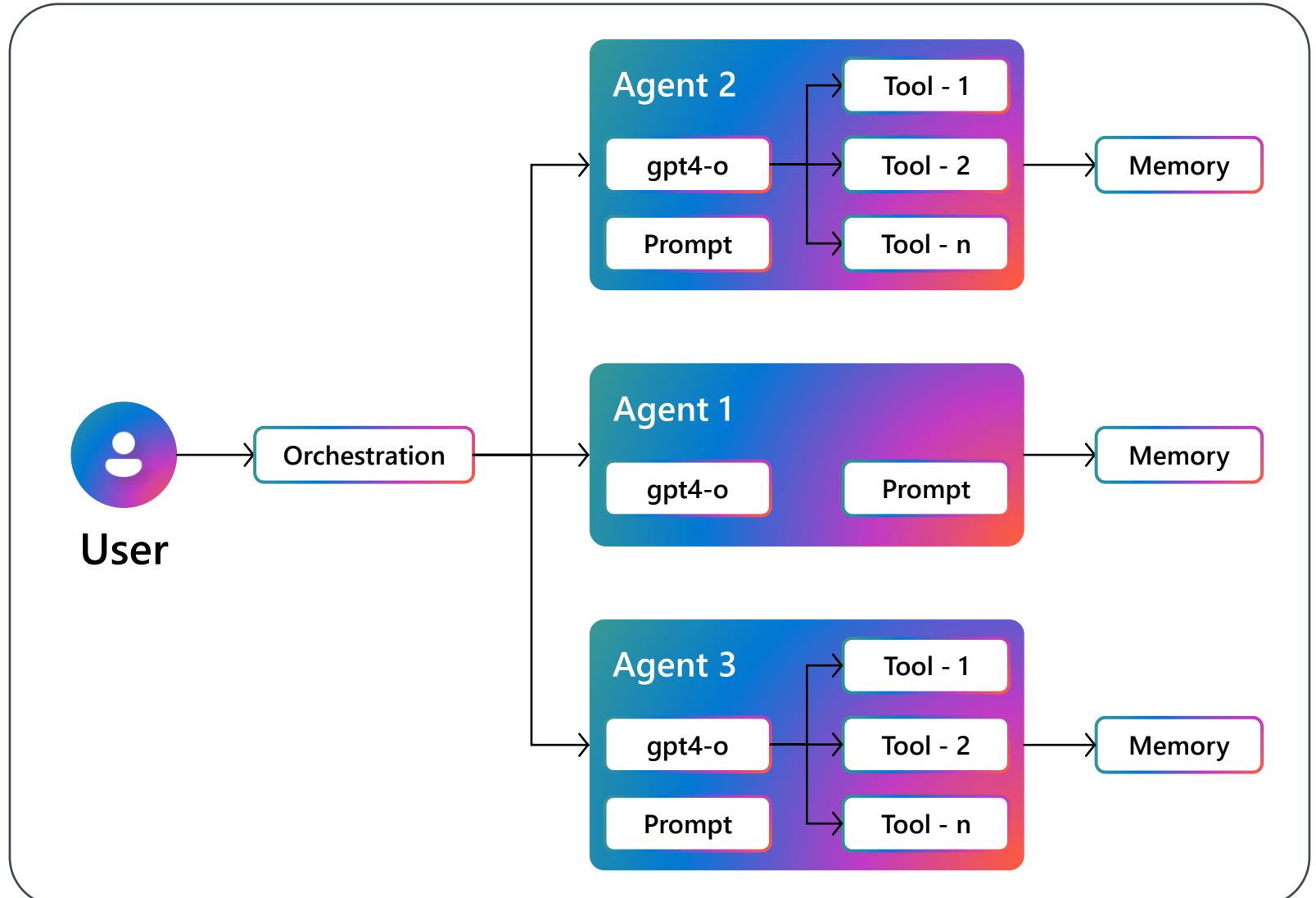
Multi Agent Logical Architecture

Each agent is specialized in different tasks or aspects of a problem

Agents can communicate and coordinate with each other. Structured orchestration is crucial

2 primary categories based on orchestration types

- Vertical Architecture
- Horizontal Architecture



Multi-Agent Collaboration

- **Specialization**
 - Different agents are configured for specific tasks. They can tackle different aspects of a complex problem.
 - Multiplying the power of a single agent.
- **More modular and easier for developers**
 - Keeping the system easy to maintain and add or remove components.
 - Increasing collaboration across different teams' Copilots.

