

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

Workshop Objectives

05

Business Scenario

06

Workshop Challenges & Team Assignments for Guided Exercises

07

Agentic Solutions

08

Demo

Workshop Results

Day 3: May 12, 2025

09

Presentation of Team Solutions

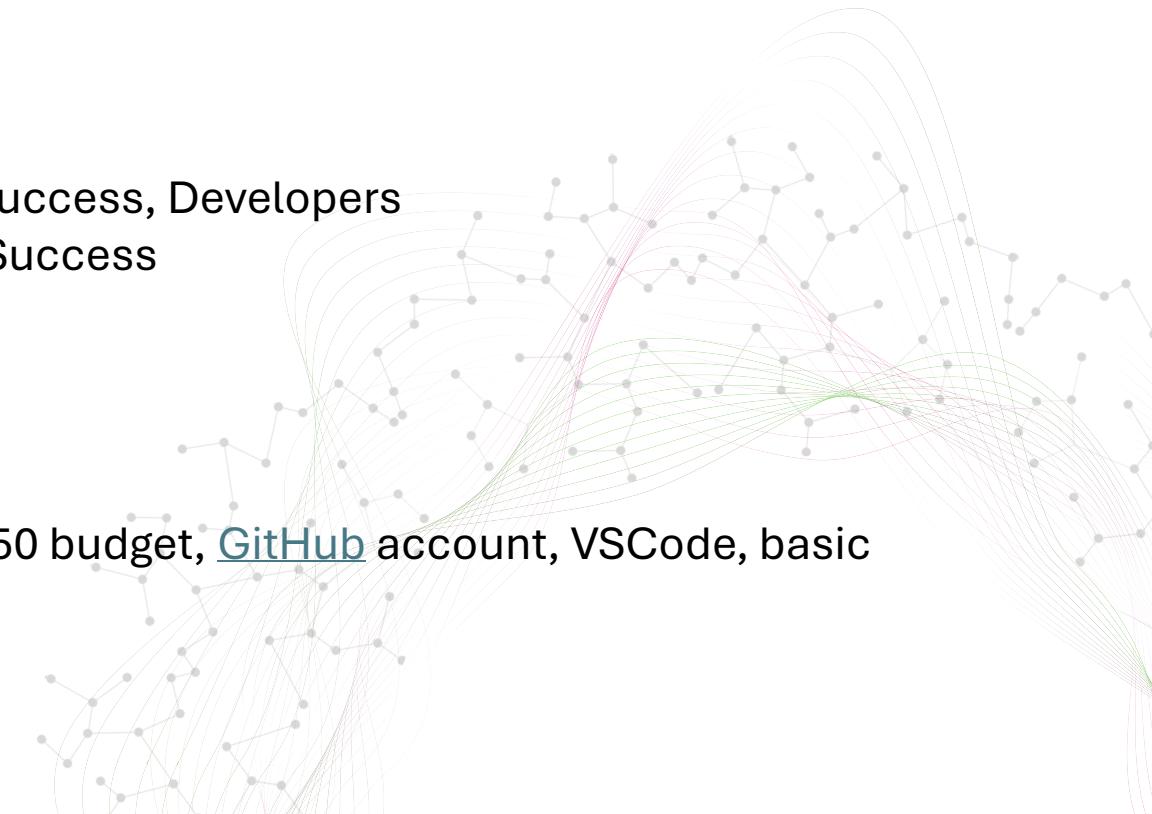
10

Wrap-up & Final Discussion



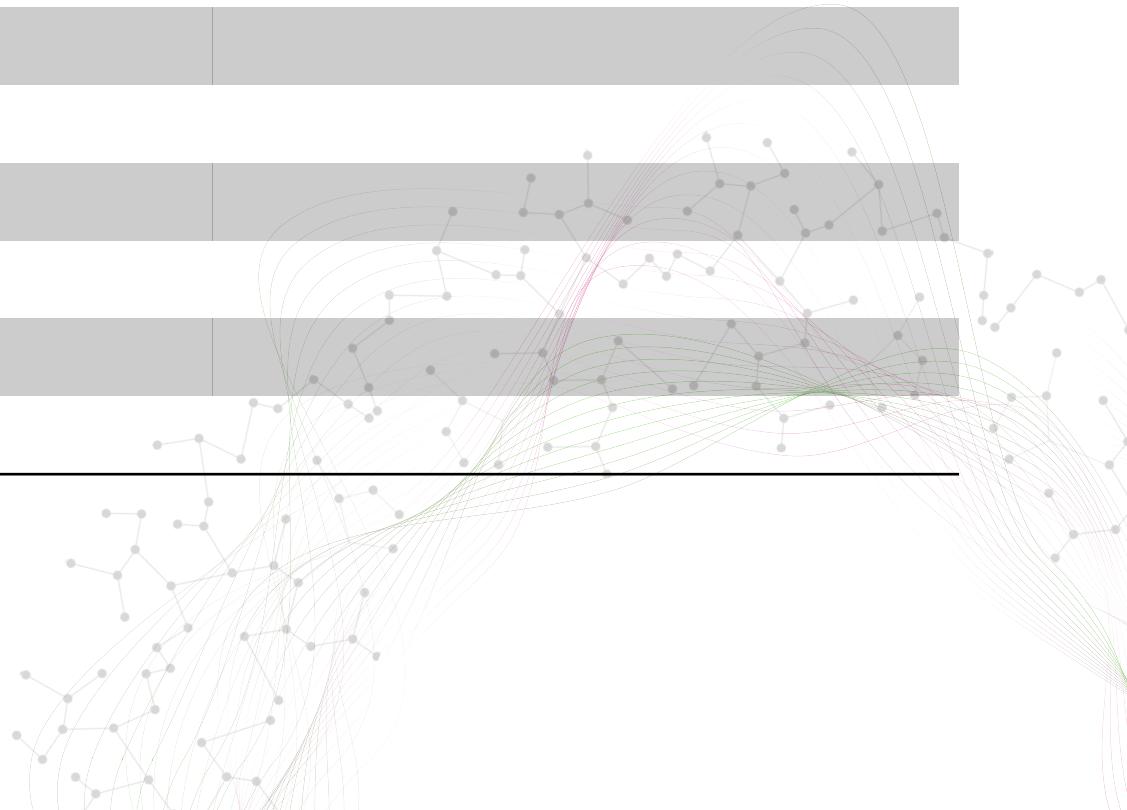
Workshop Overview

- What you will learn
 - By the end of the workshop, you will learn
 - How to build an AI Agent and explore its [tools](#)
 - What are different AI Agent Platforms and how to use them
 - Hands-on experience building single and multi-agent systems
- Getting Started
 - [README](#)
- Target Audience
 - Day1: Executives, Product Management, Customer Success, Developers
 - Day2: Developers, Product Management, Customer Success
 - Day3: Developers
- Pre-requisites
 - Day1: Basic understanding of AI and LLMs
 - Day2 & Day3: Access to an Azure subscription with \$50 budget, [GitHub](#) account, VSCode, basic familiarity with python



Workshop Support Team

Name	Contact	Notes
James Nguyen		
Anil Dwarkanath		
Nicole Serafino		
Patrick O'Malley		
Heena Ugale		
Aditya Agrawal		
Claire Rehfuss		
Kirby Repko		





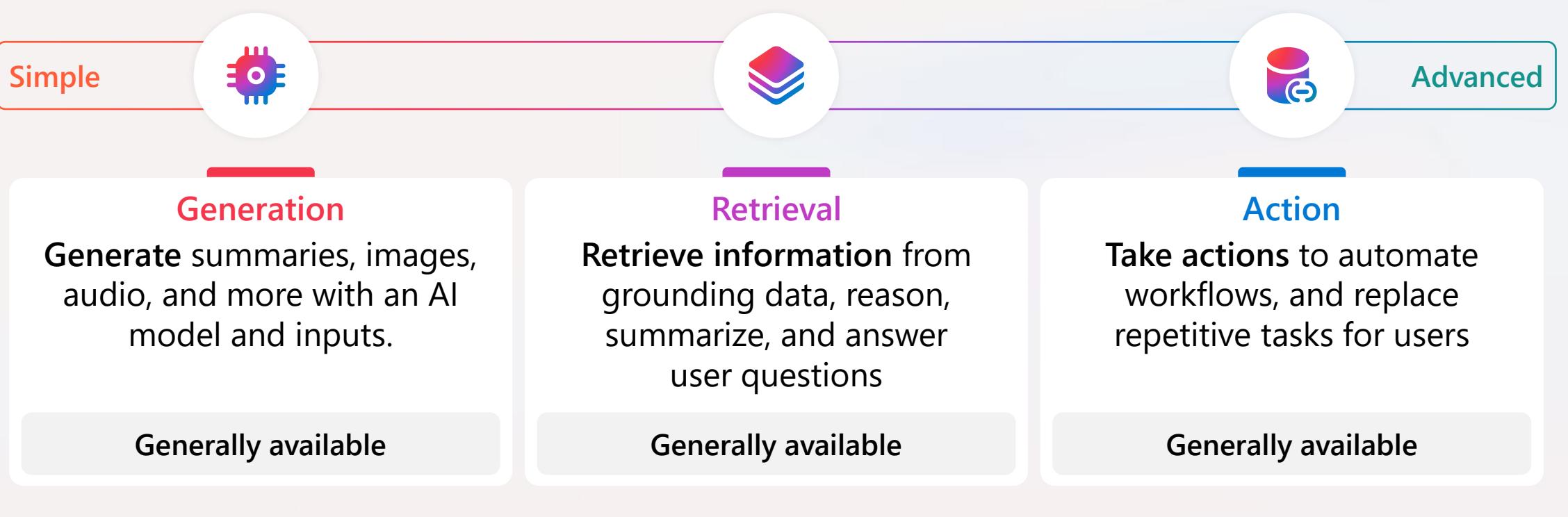
Day 1

Introduction to Agentic AI

What are agents?

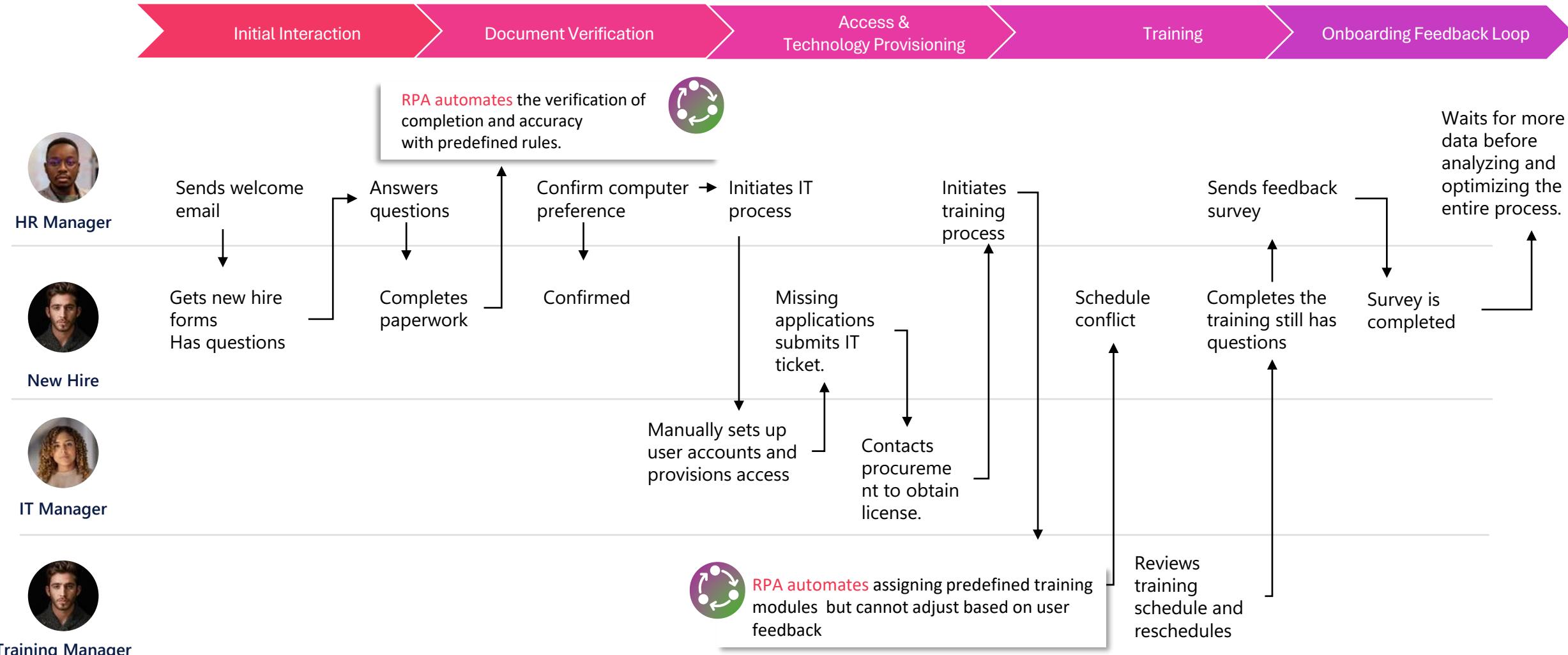
AI designed to perform a task

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



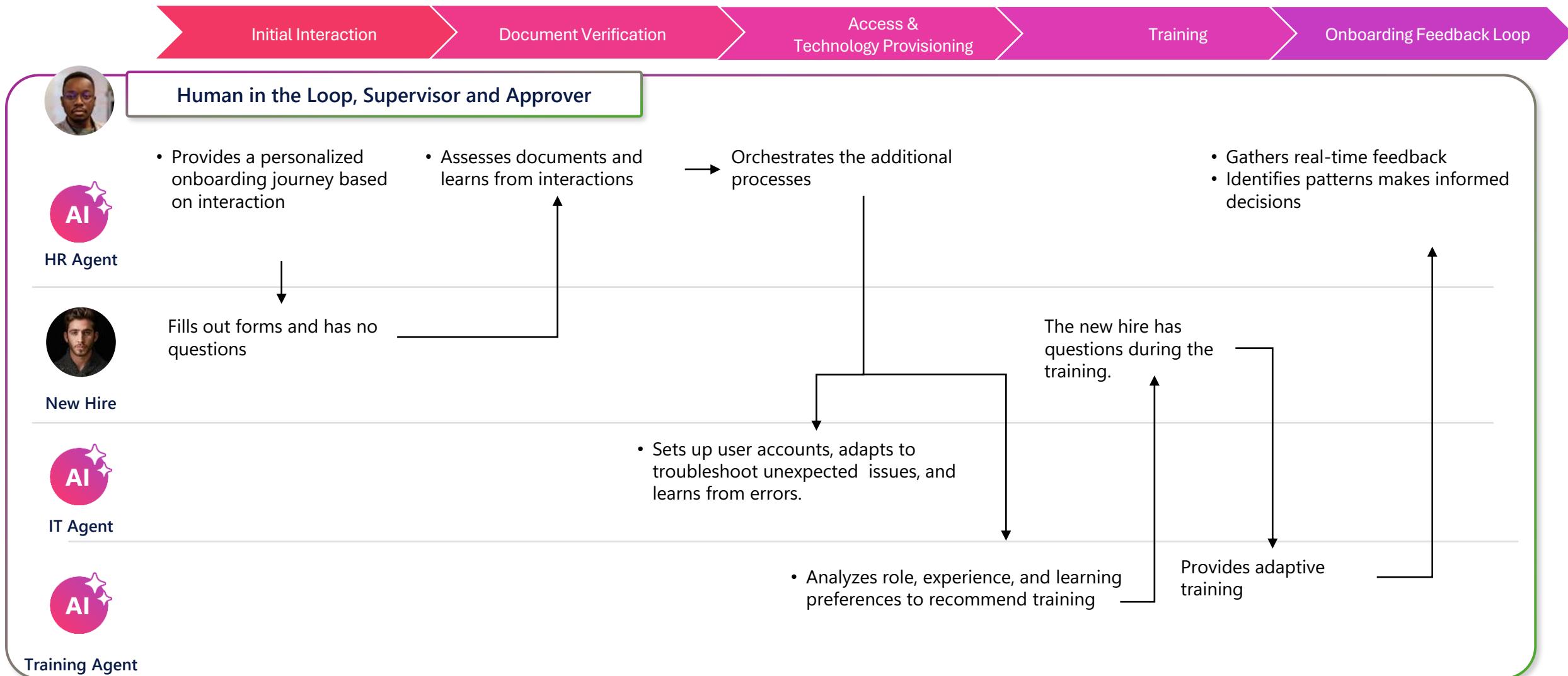
TODAY – Workflow Automation

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

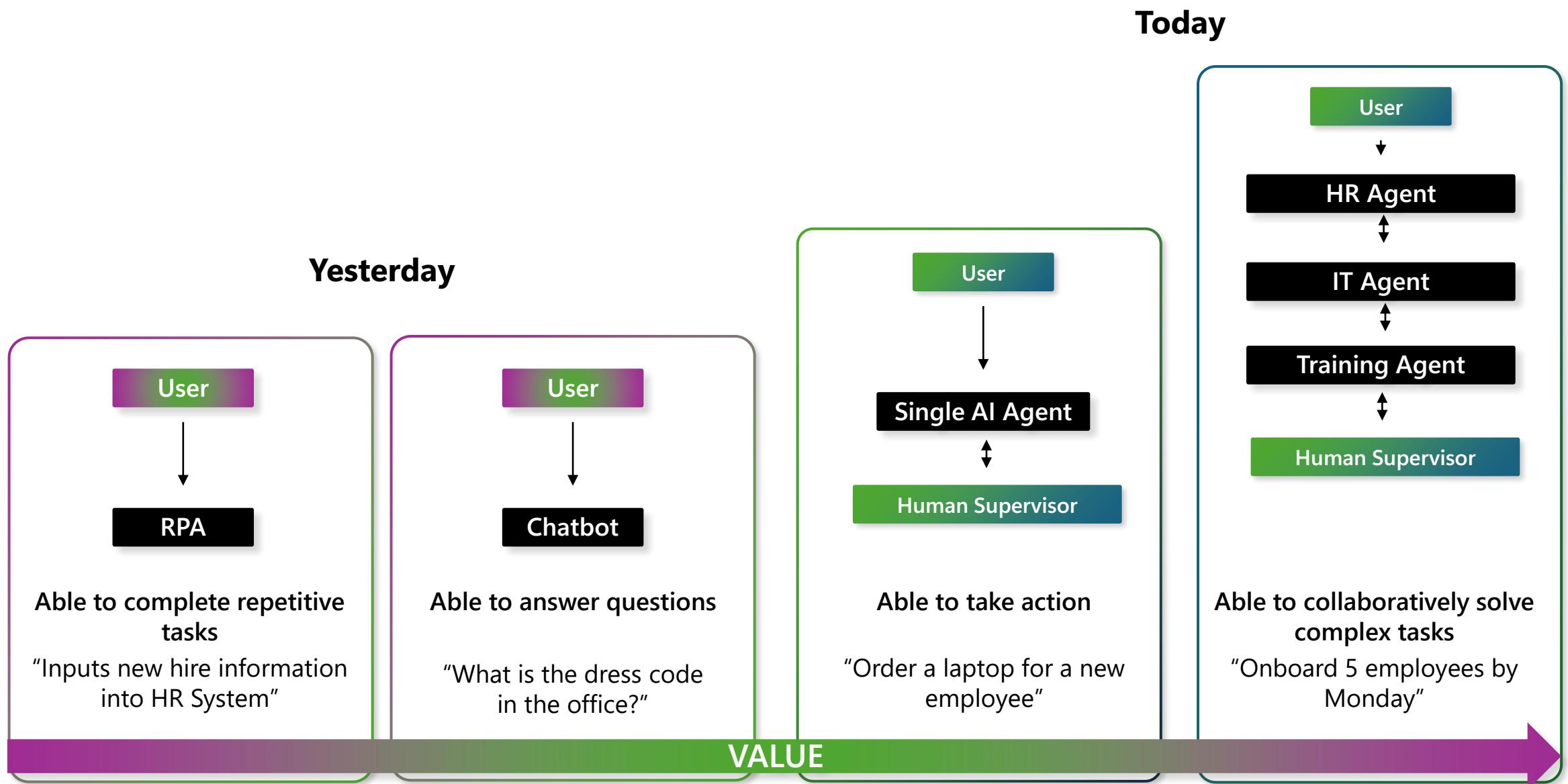


TOMORROW – AI Agents

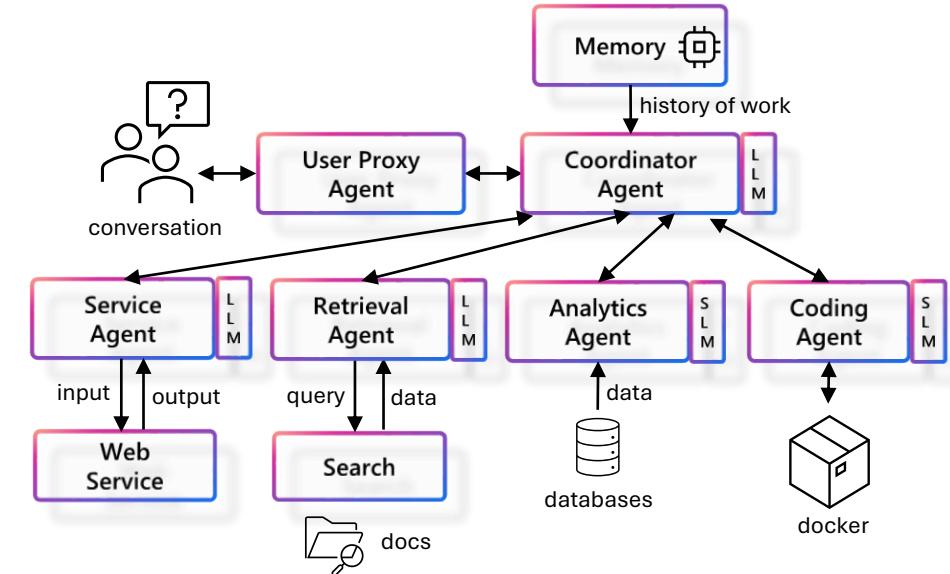
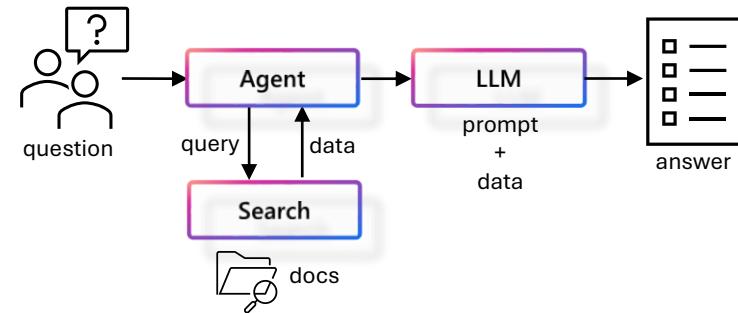
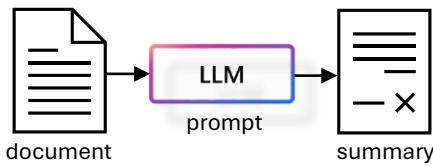
With AI Agents, these steps can be fully automated for the first time.



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

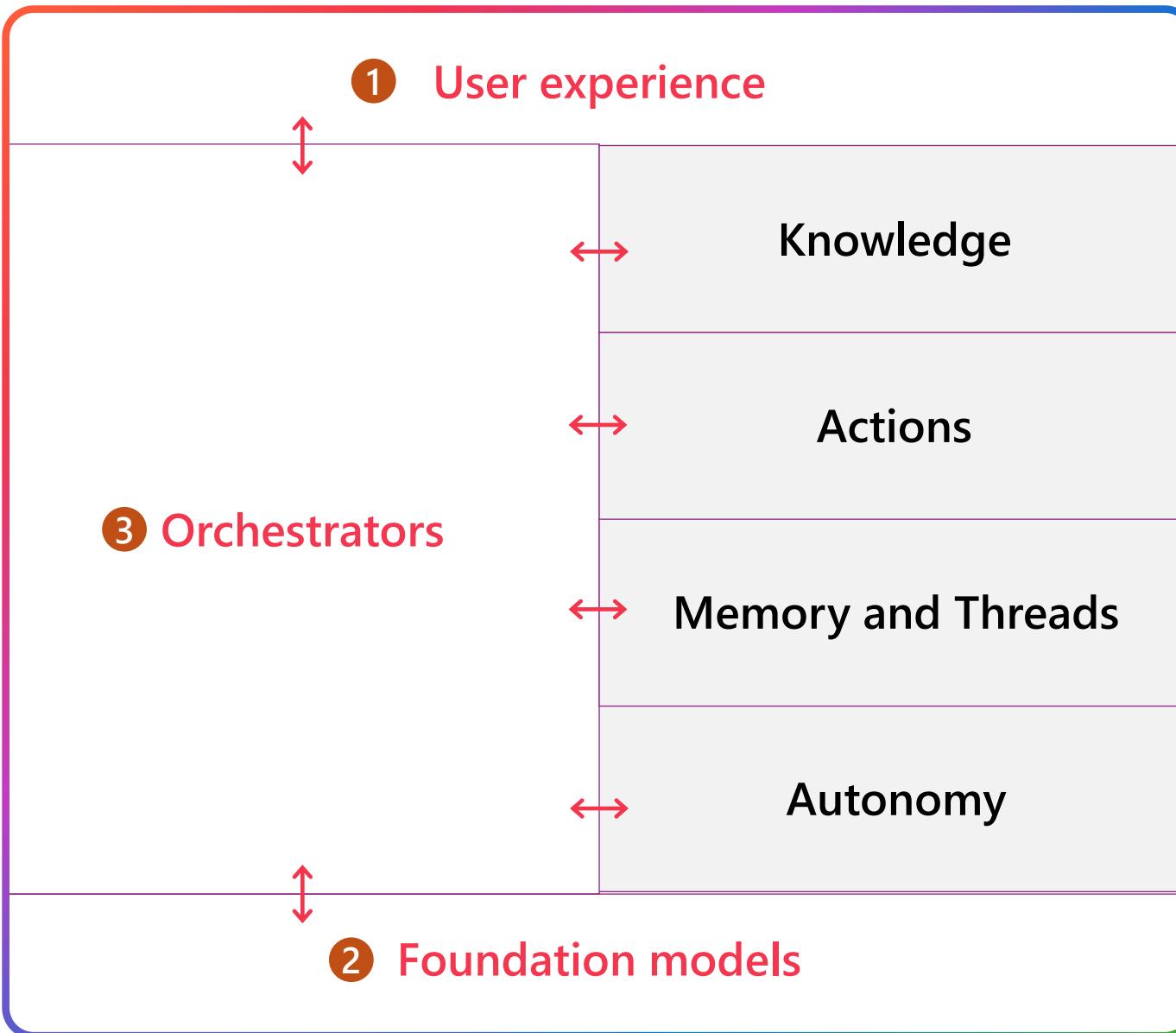
Multi-agent Systems

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

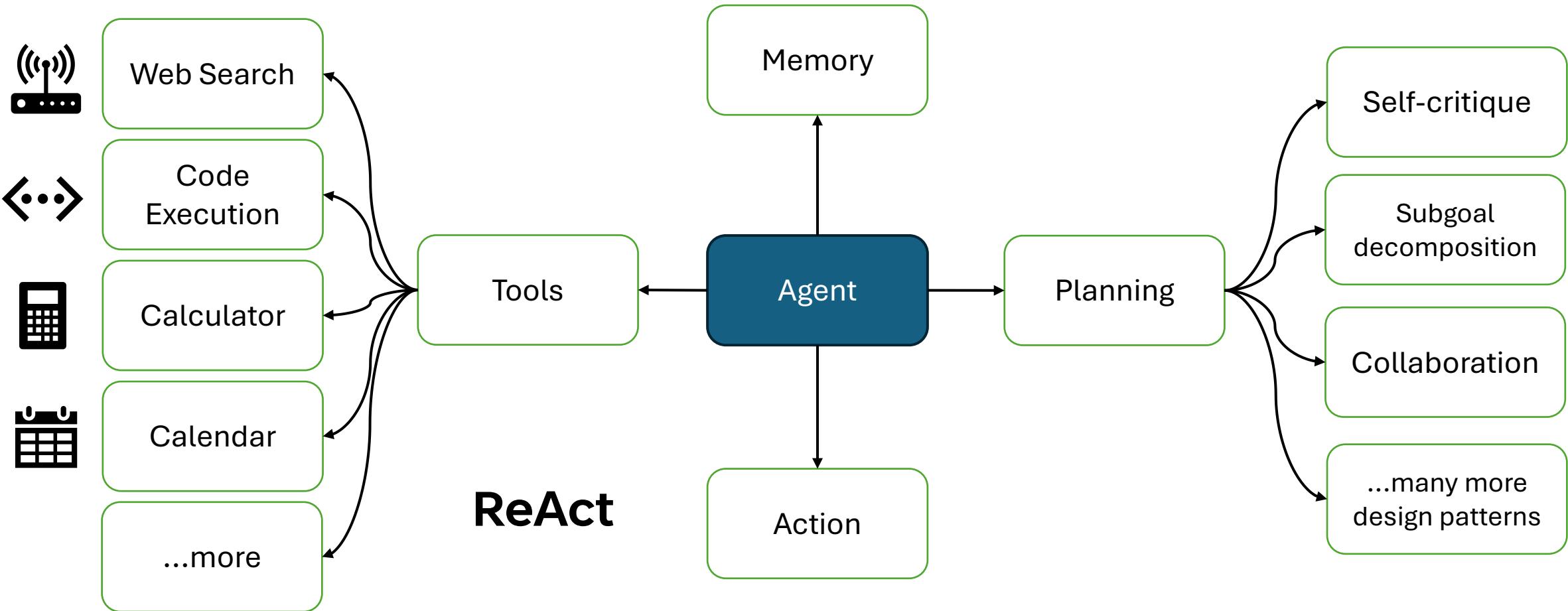
VALUE

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.

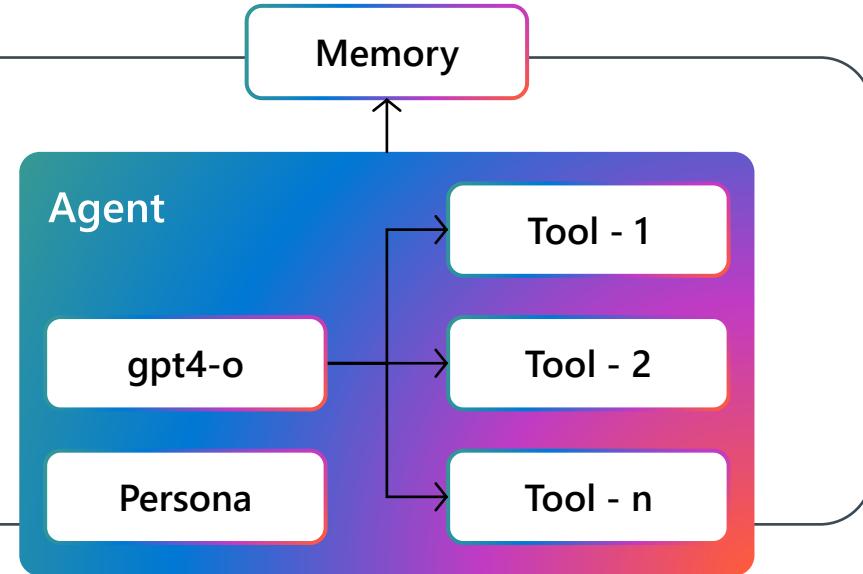
Agentic AI capabilities



Agent Abstractions - Agent First-Class Citizen

Agent as high-level abstraction

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

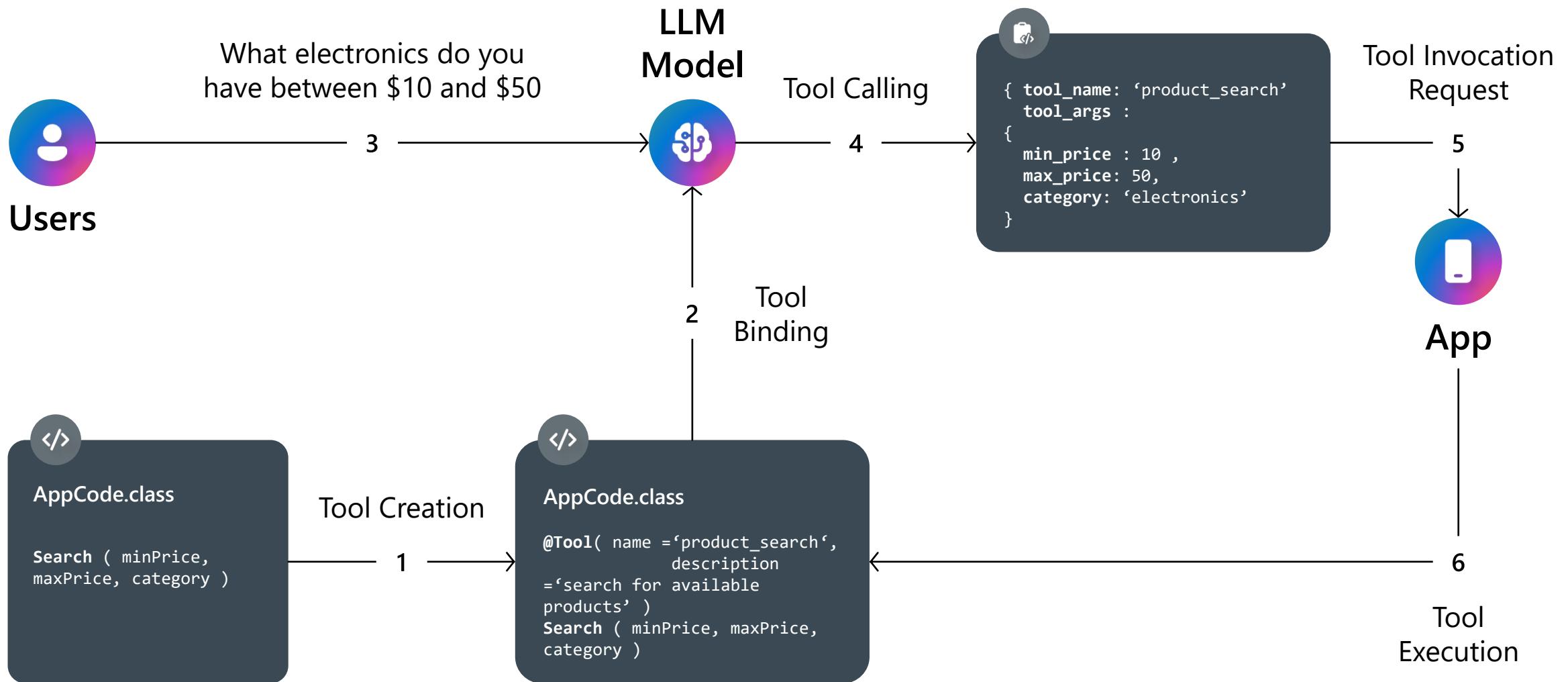


Agent Chat as layer for collaboration

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

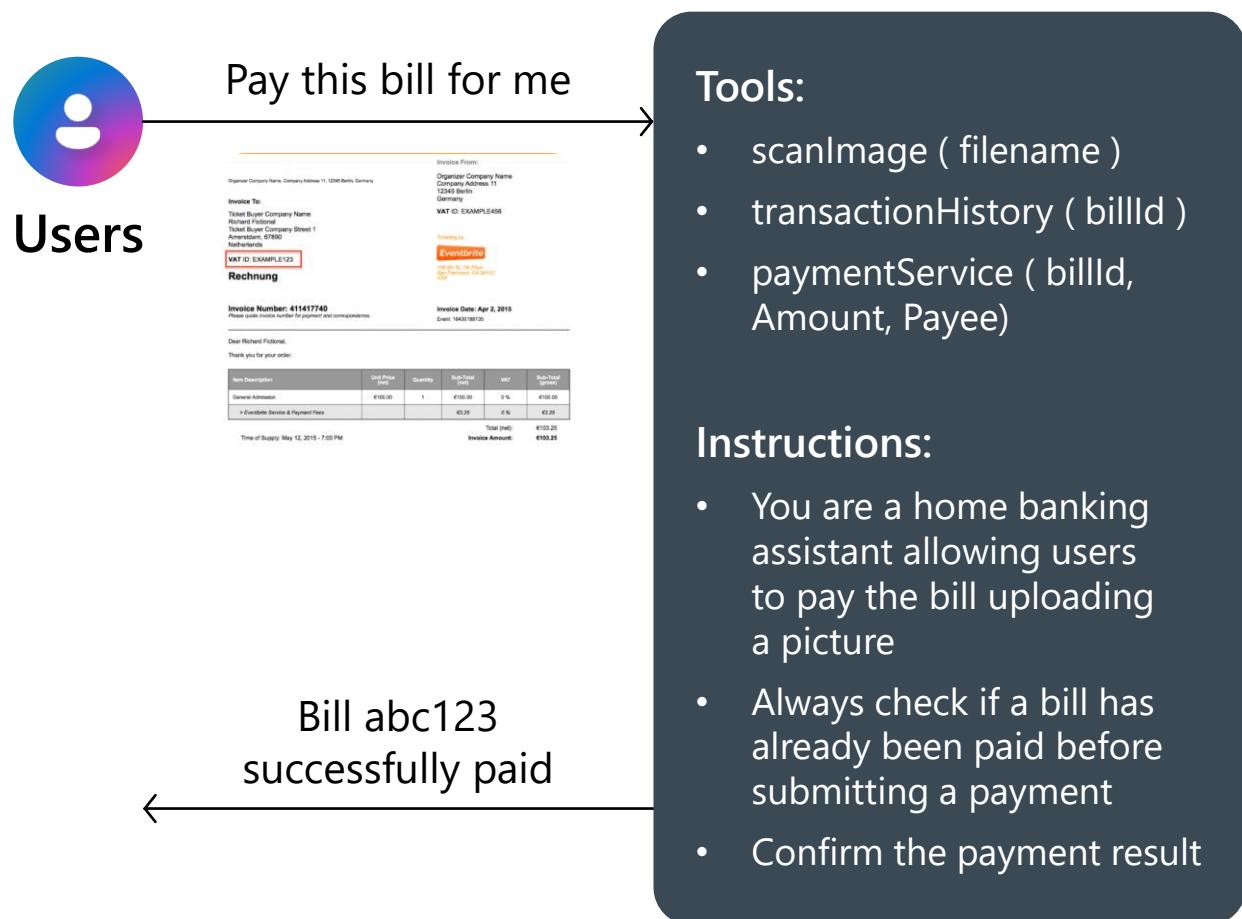


Agentic Pattern - Tools Calling

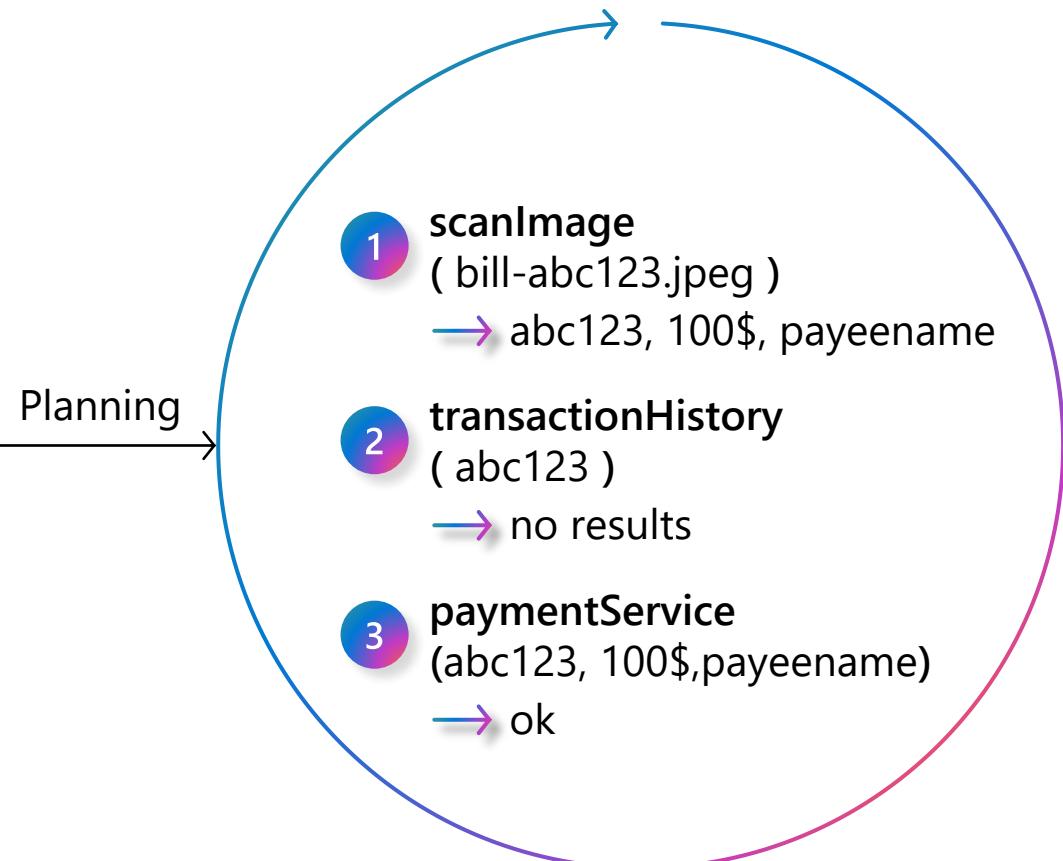


Agentic Pattern - ReAct Planning with Tools

Calling Payment Agent



`while (new tools execution request)`



Agentic Pattern - Memory

Short Term

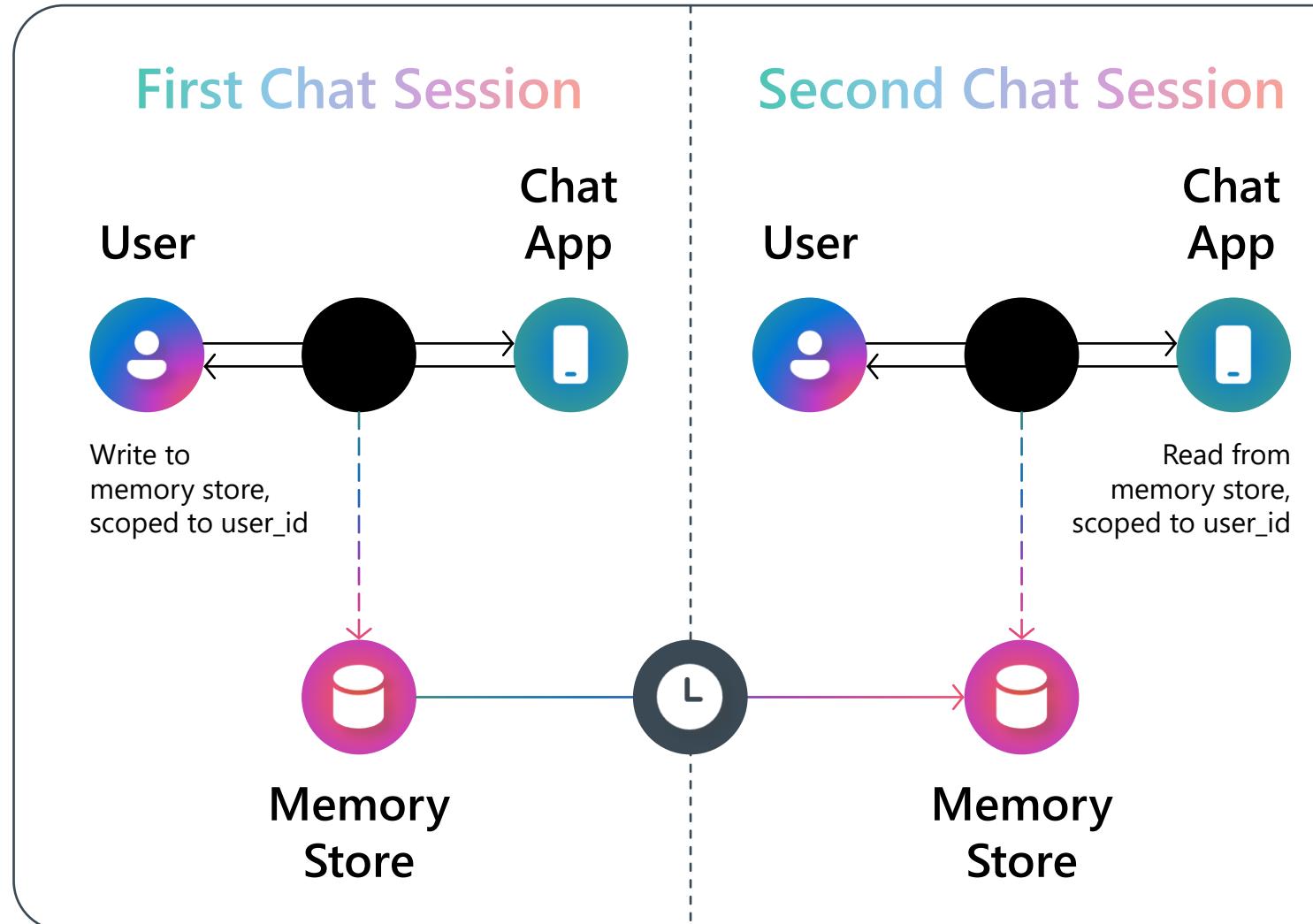
- Access steps info in one loop iteration
- Shared state context
- Chat history

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)



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

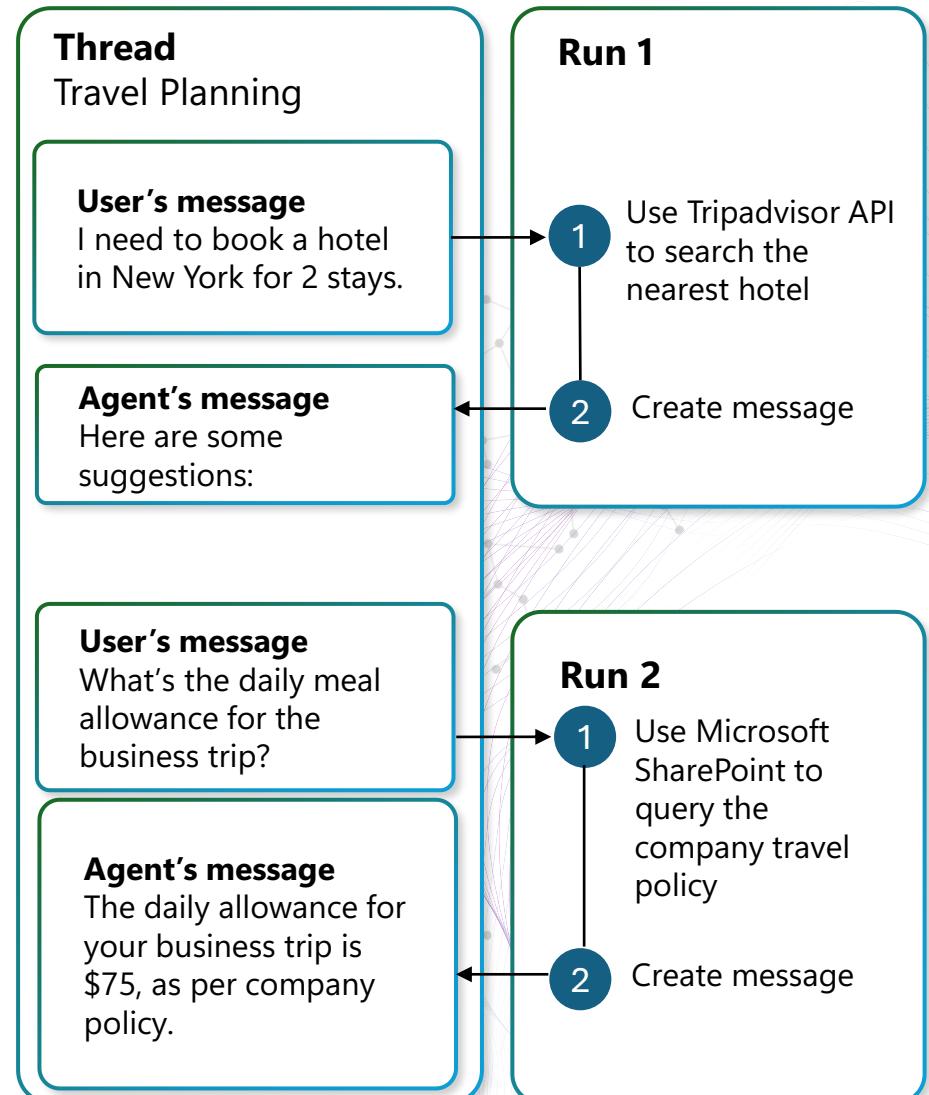


Thread – Providing History to Conversations

Threads store conversation history. Currently, agents store threads in a managed, multi-tenant Cosmos DB. Soon, customers will be able to bring their own (BYO) Cosmos DB as a part of the standard setup, ensuring all thread messages are stored in their customer tenant for enhanced security.

Definitions

- Thread – a conversation session between an agent and a user. Threads store Messages and automatically handle truncation to fit content into a model's context.
- BYO Cosmos DB – an enterprise feature that provides secure, single-tenant data control, allowing you to maintain complete control over your customer data.



Thread Storage VS Memory

	Thread 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

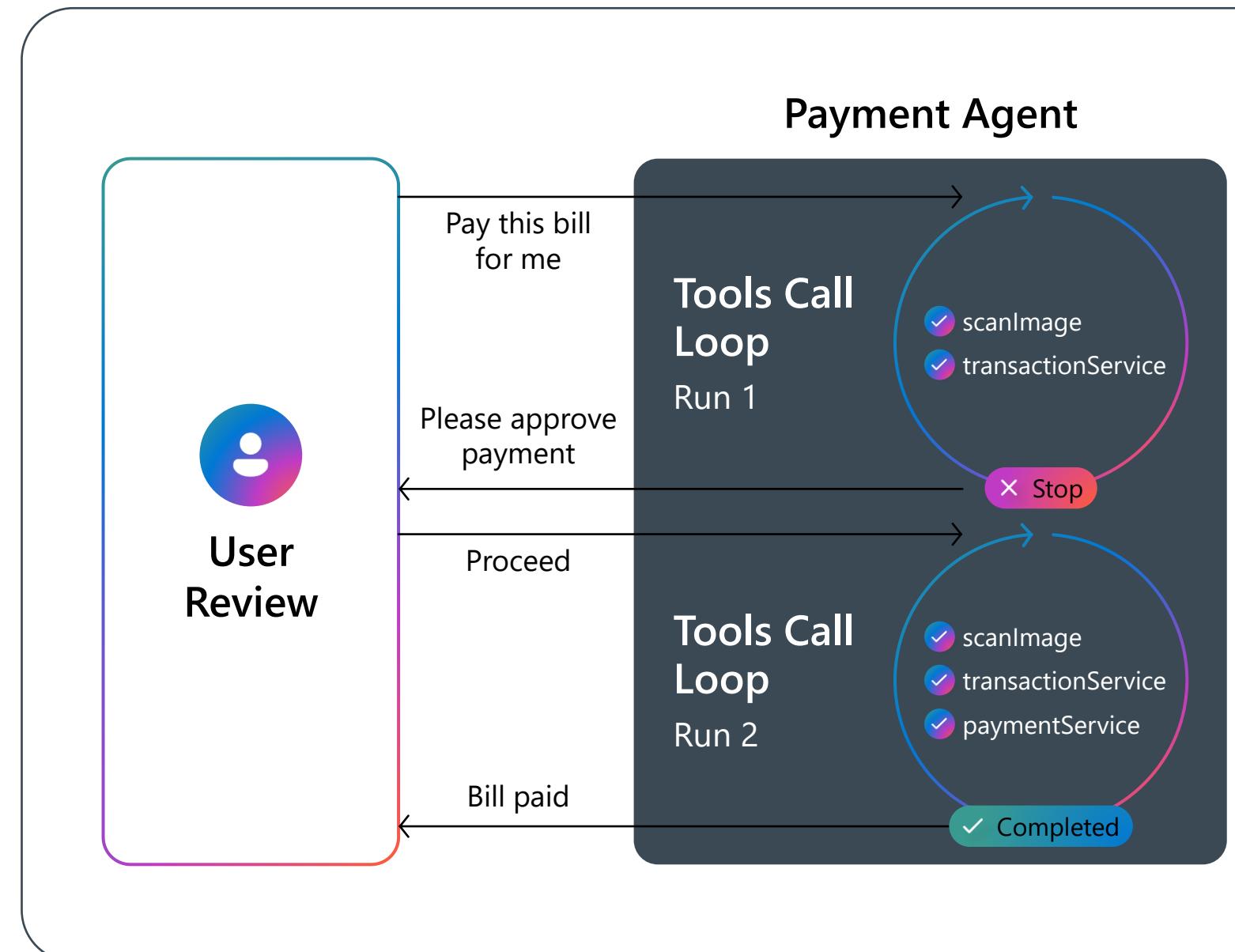
Agentic Pattern - Flow control

Looping Termination

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

Human in the loop

- Action execution approval
- Escalation
- Data review





Day 1

Agentic System Design Patterns

Single Agent

Multi-domain

Multi-agent

Design Patterns Explored in Workshop

Intelligent Single Agent



Multi-domain Multi-agent

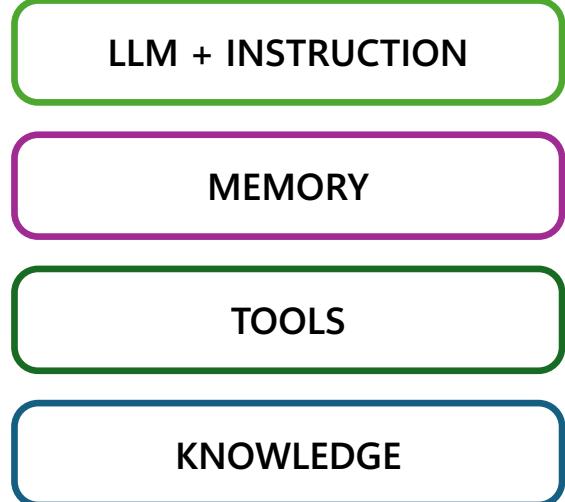


Collaborative Multi-Agent

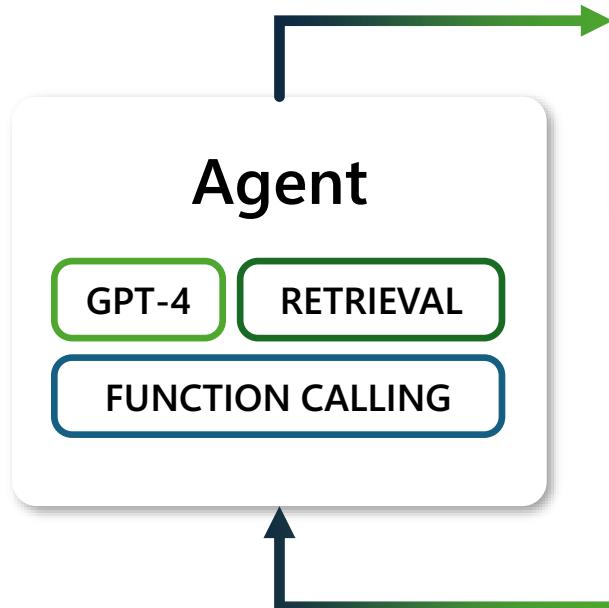


Single vs Multi-agent

Single Agent



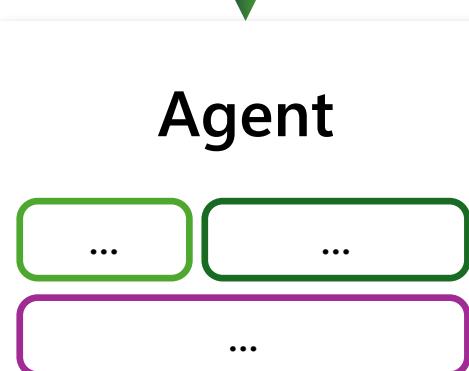
Multi-Agent



Agent

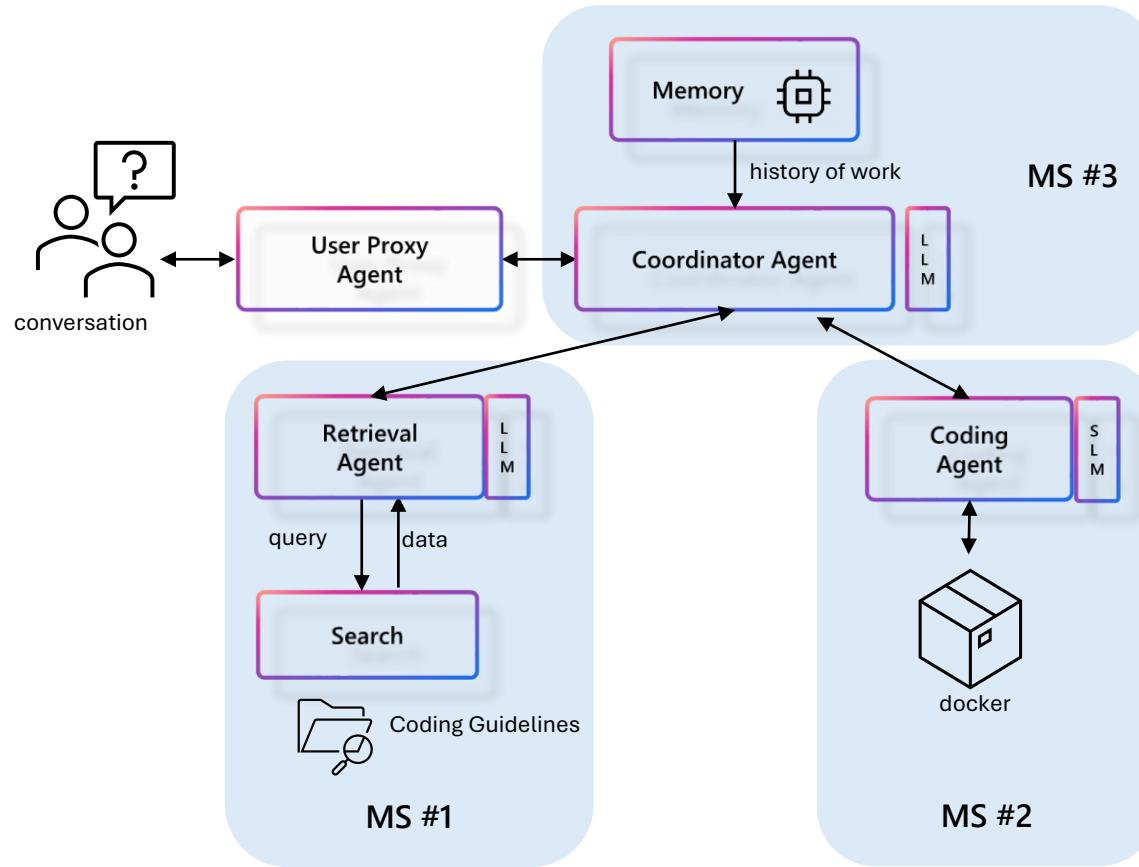


Agent



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

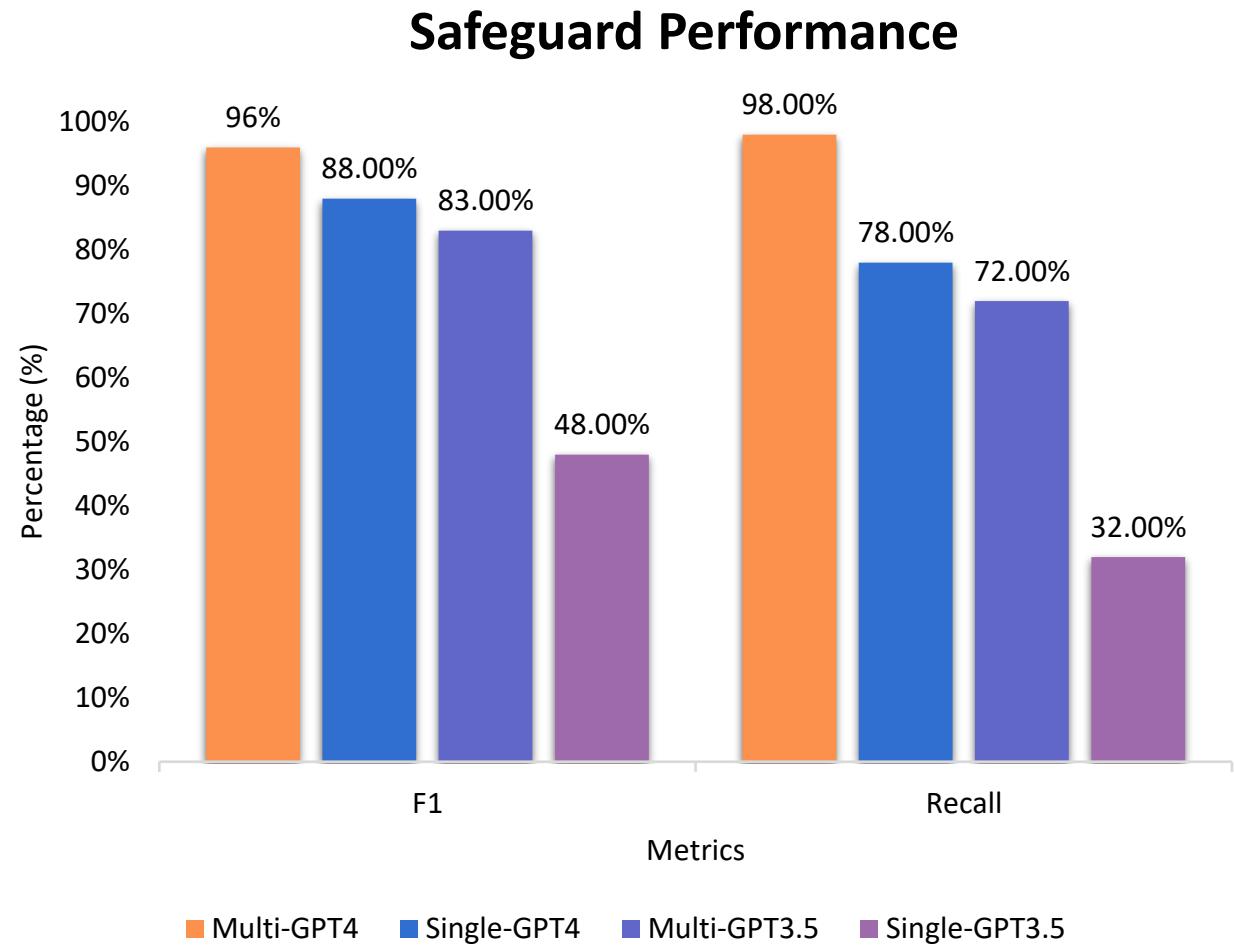
- ✓ Adaptive planning within scope of existing tightly scoped skills (agents)
- ✓ Handles ambiguity by discussing and refining requirements with human
- ✓ Memory to handle complex long running execution of a plan
- ✓ Effective inter agent communications
- ✓ Test, monitor, release & maintain each agent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

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

Multi agent systems can solve more complex tasks

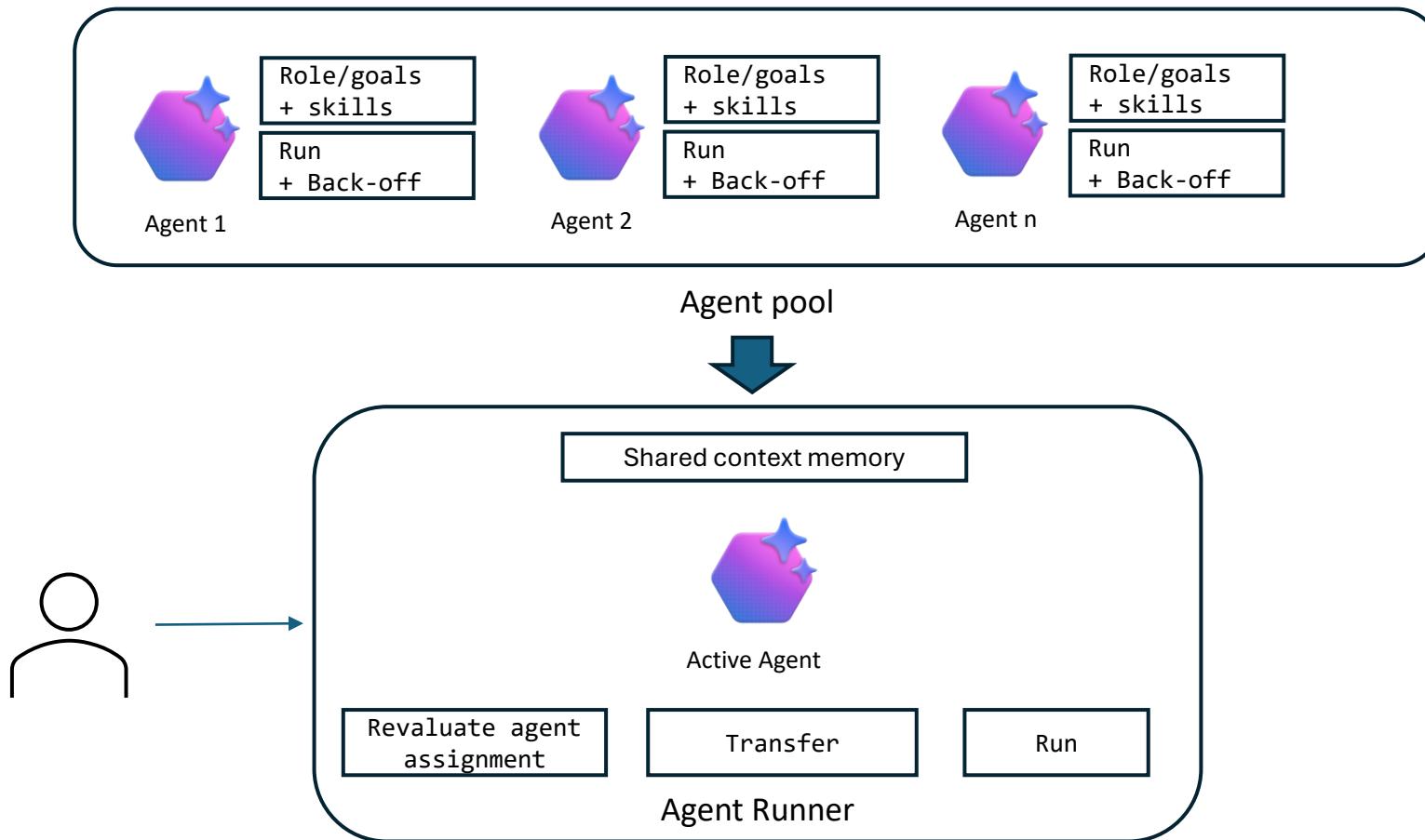
Single agents can solve a wide range of problems and are simpler to implement.

- 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



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

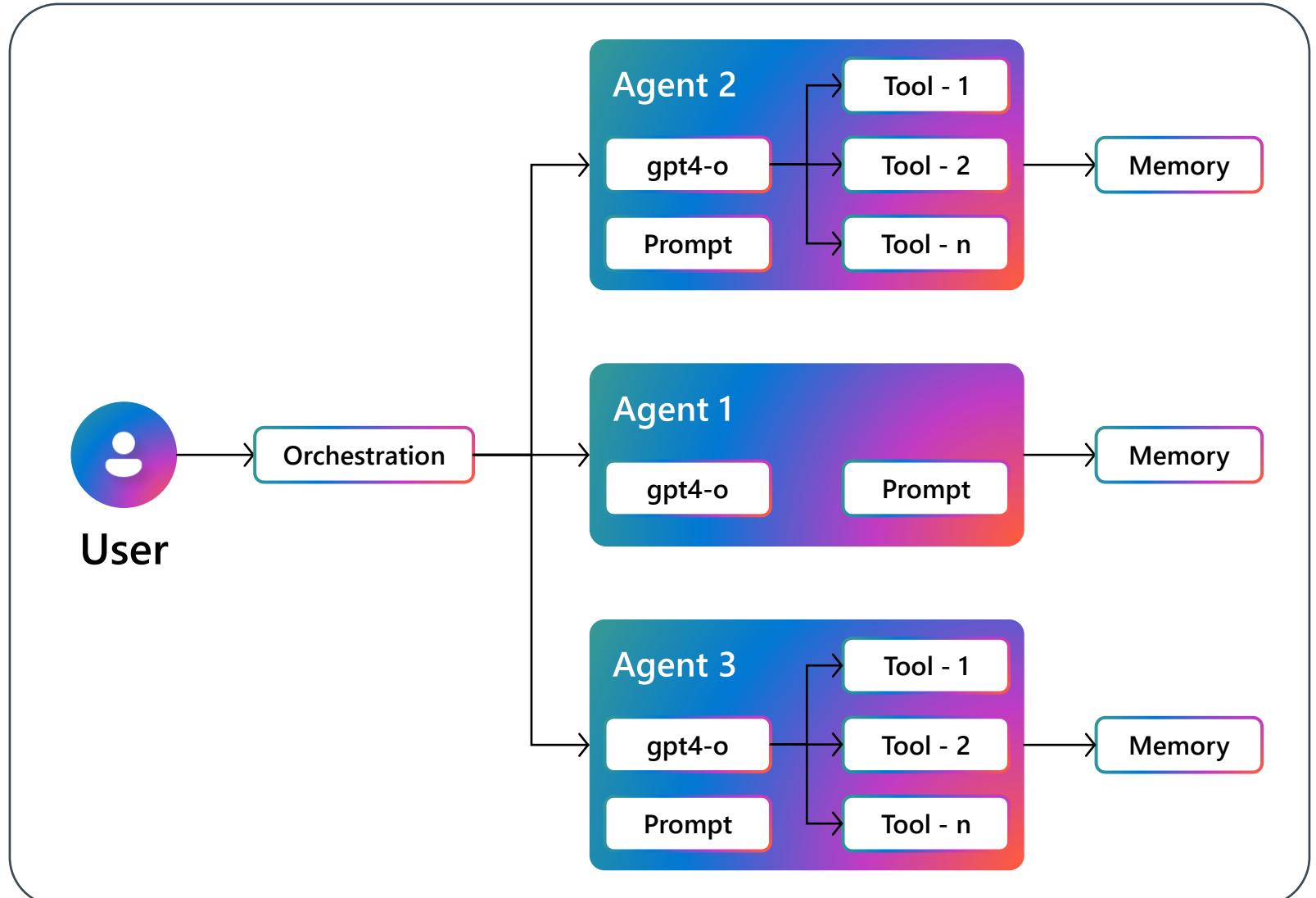
Multi Agent Logical Architecture

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

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

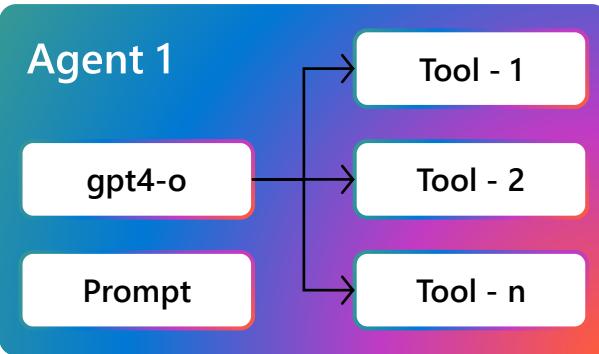
2 primary categories based on orchestration types

- Vertical Architecture
- Horizontal Architecture

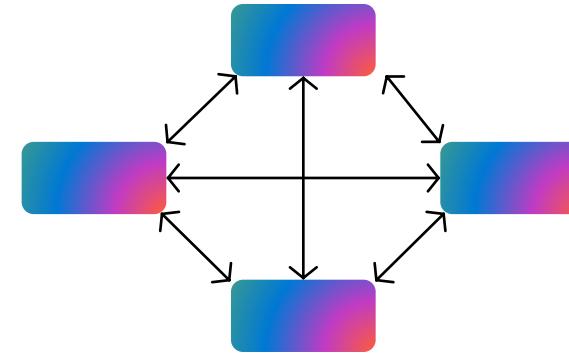


Agents orchestration and communication styles

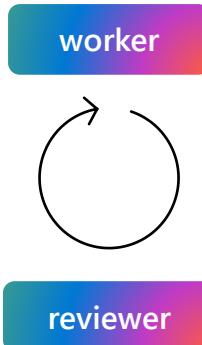
Single Agent



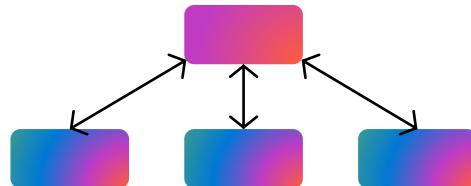
Network



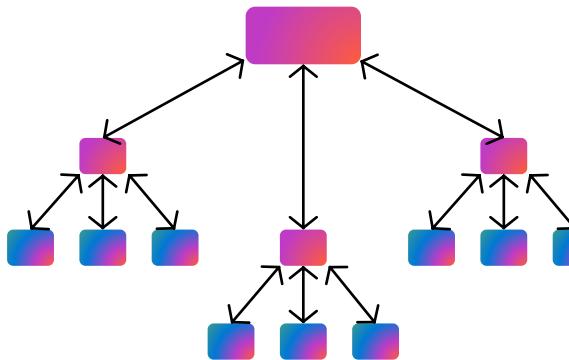
Reflection



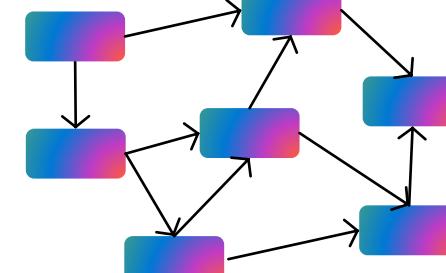
Supervisor



Hierarchical

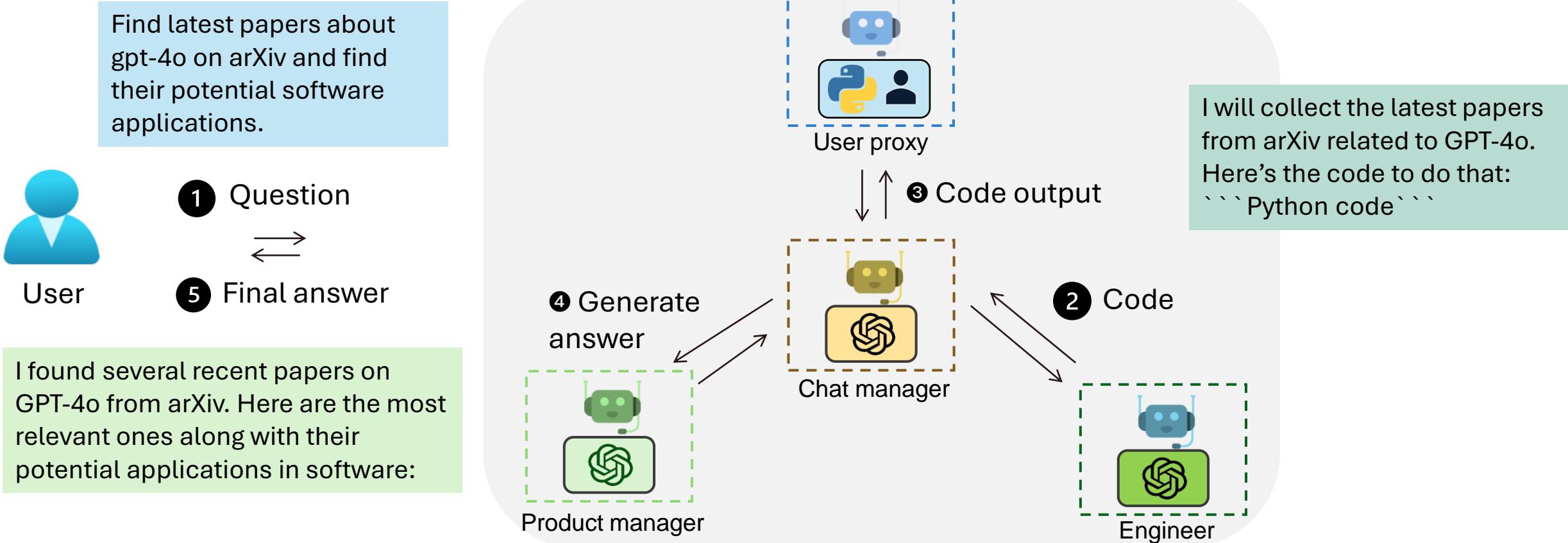


Custom



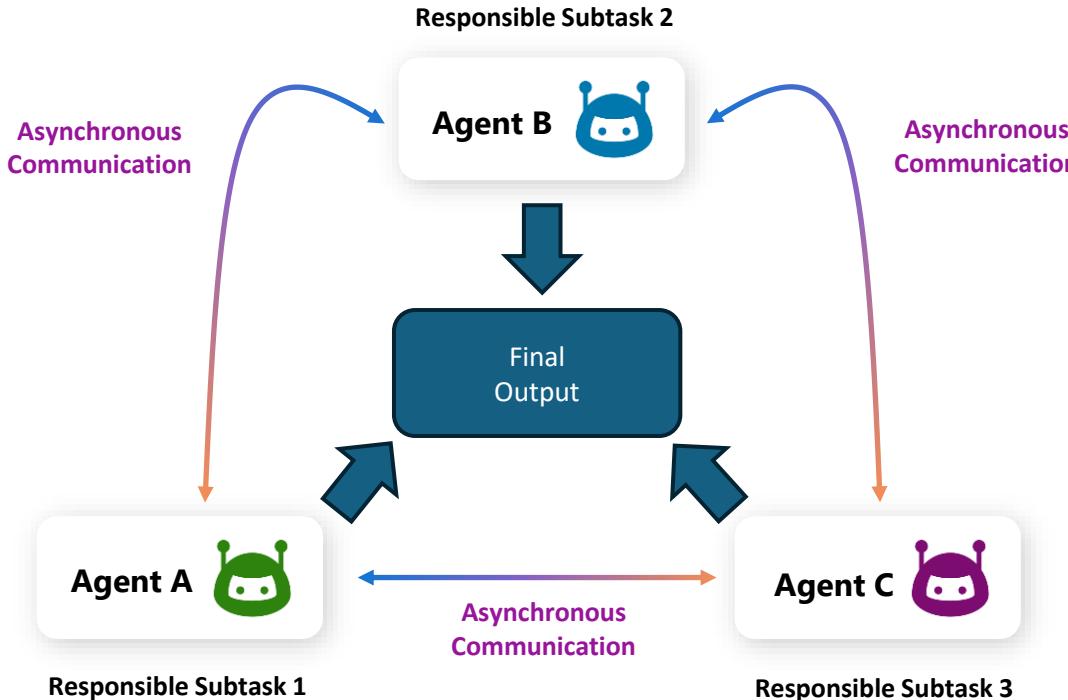
Multi-Agent Collaboration

- **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.



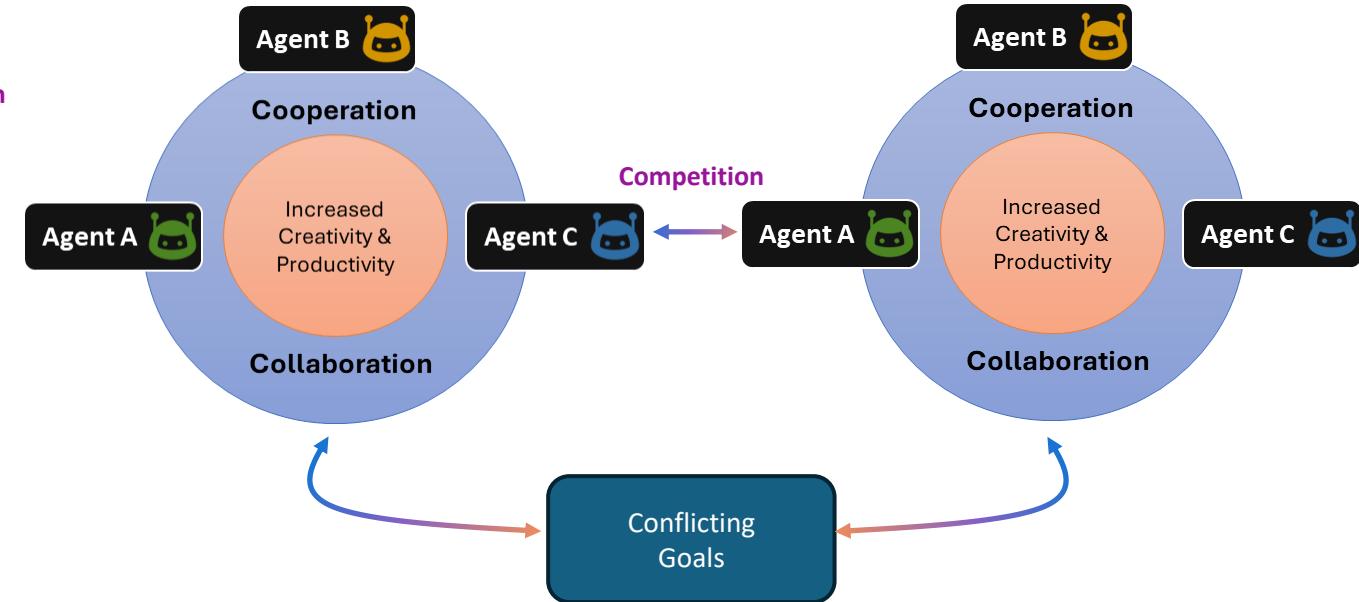
Cooperative vs Competitive Agents

Cooperative based-learning



Competition based-learning

Seems not useful in a machine scenario



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

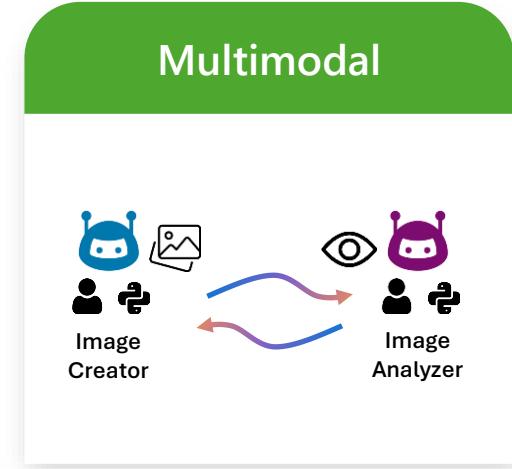
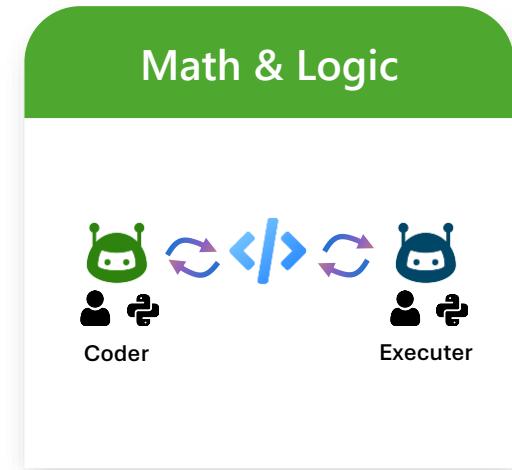
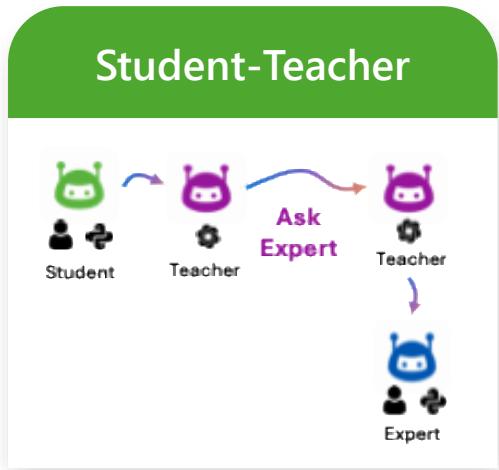
In contrast, competitive agents operate with individual goals in adversarial settings, focusing on optimizing their own outcomes, often with strategic communication and planning.

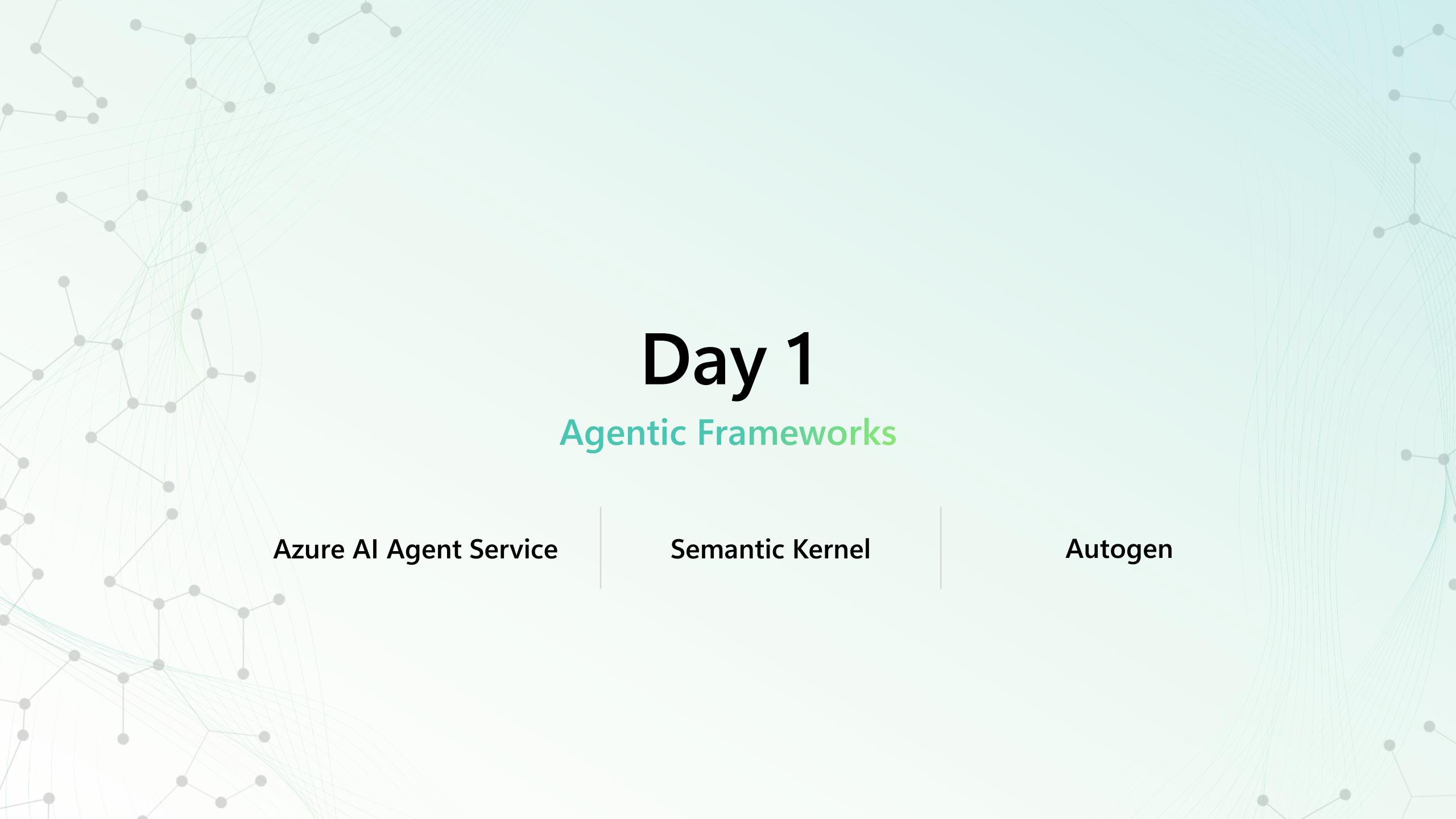
Resource Allocation

Decision Making

Realism

Multi Agent Examples



A faint, light-gray network graph with numerous nodes and edges, forming a complex web-like pattern across the entire slide.

Day 1

Agentic Frameworks

Azure AI Agent Service

Semantic Kernel

Autogen

Role of Frameworks

Abstraction Layers

Reduce Complexity
of integrating tools
and workflows

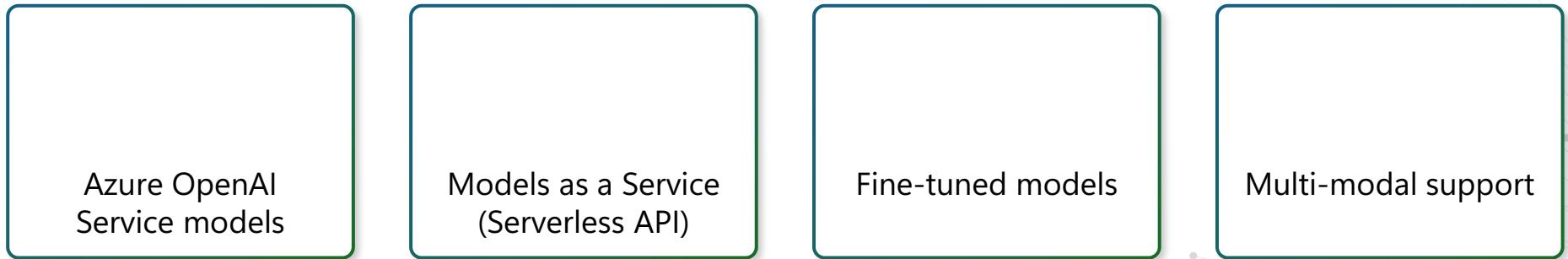
Built in Functionality

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 (x) Collections (x) Industry (x) Deployment options (x) Inference tasks (x) Fine-tuning tasks (x) Licenses (x)

Search

Models 1795

gpt-4o-realtime-preview	openai-whisper-large-v3	openai-whisper-large	gpt-4
Chat completion	Speech recognition	Speech recognition	Chat completion
gpt-35-turbo	o1-preview	o1-mini	gpt-4o-mini
Chat completion	Chat completion	Chat completion	Chat completion
gpt-4o	gpt-4-32k	gpt-35-turbo-instruct	gpt-35-turbo-16k
Chat completion	Chat completion	Chat completion	Chat completion
dall-e-3	dall-e-2	whisper	tts-hd
Text to image	Text to image	Speech recognition	Text to speech
tts	text-embedding-3-small	text-embedding-3-large	Phi-3.5-mini-instruct
Text to speech	Embeddings	Embeddings	Chat completion
Phi-3.5-MoE-instruct	Phi-3-mini-4k-instruct	Phi-3-medium-4k-instruct	Phi-3-mini-128k-instruct
Chat completion	Chat completion	Chat completion	Chat completion
Phi-3-medium-128k-instruct	Phi-3-small-8k-instruct	Phi-3-small-128k-instruct	Phi-3.5-vision-Instruct
Chat completion	Chat completion	Chat completion	Chat completion

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](#)

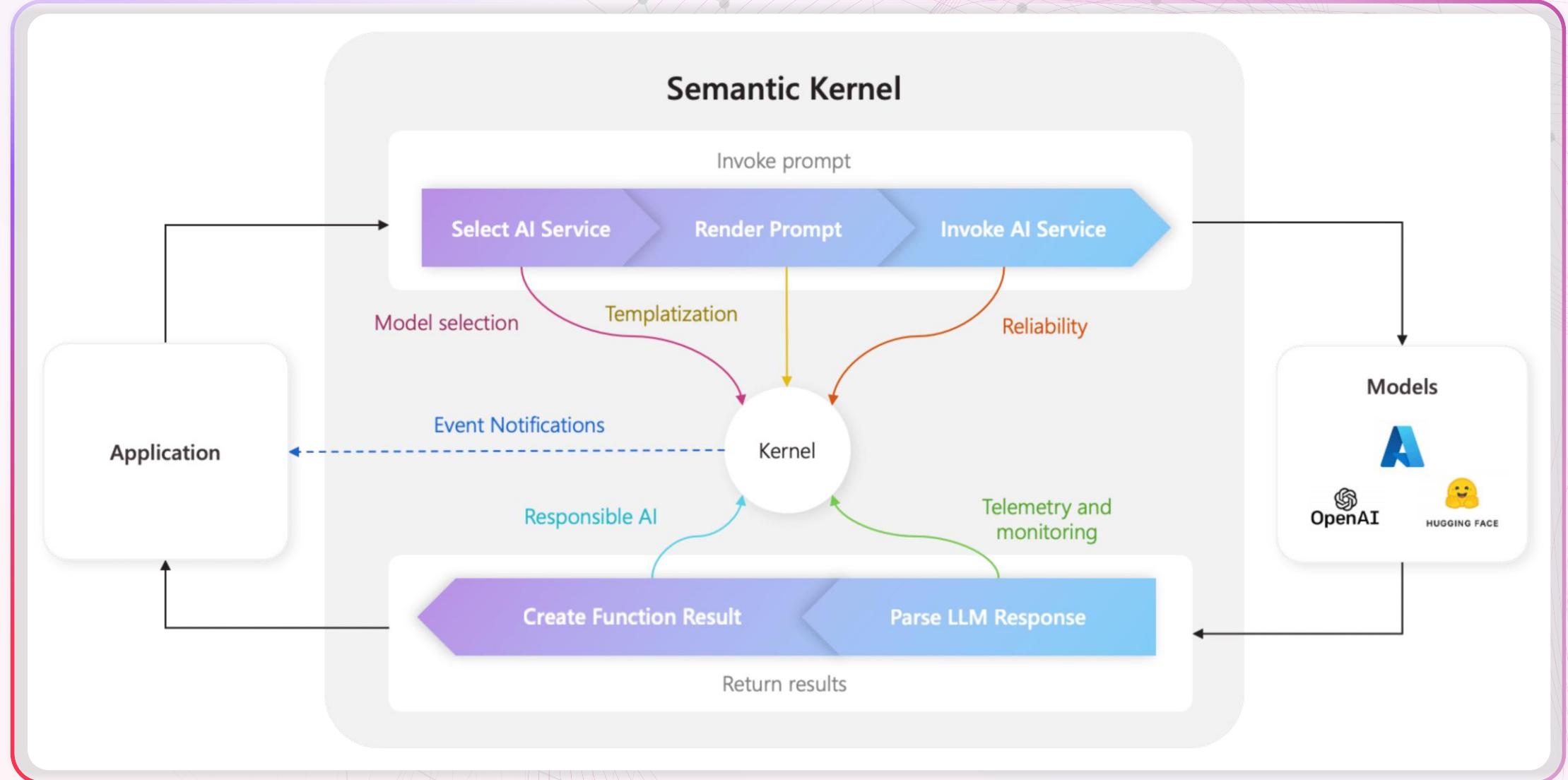


Semantic Kernel

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

The kernel is at the center



Semantic Kernel Agents

Reflection multi-agent

- Coding agents
- Digital content creation agents

ChatCompletionAgent

- Name
- Instructions
- Kernel (Services, Plugins)

Reflection loop with 2 agents

- Blog post writer
- Blog post reviewer

```
from agents.base_agent import BaseAgent
from semantic_kernel.agents import ChatCompletionAgent, ChatHistoryAgentThread
from semantic_kernel.connectors.ai.open_ai import AzureChatCompletion
from semantic_kernel.connectors.mcp import MCPSsePlugin

async def chat_async(self, prompt: str) -> str:
    # Ensure agent/tools are ready and process the prompt.
    await self._setup_agent()

    response = await self._agent.get_response(messages=prompt, thread=self._thread)
    response_content = str(response.content)

    self._thread = response.thread
    if self._thread:
        self._setstate({"thread": self._thread})

    messages = [
        {"role": "user", "content": prompt},
        {"role": "assistant", "content": response_content},
    ]
    self.append_to_chat_history(messages)

    return response_content
```



Autogen

- Powerful multi-agent framework with prebuild 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.

