

# Mastering Azure AI Foundry with Agents



**Lino Tadros, RD**  
Co-Founder & CEO - Tahubu





Lino Tadros

Co-founder & Chief Executive Officer

 lino@tahubu.com

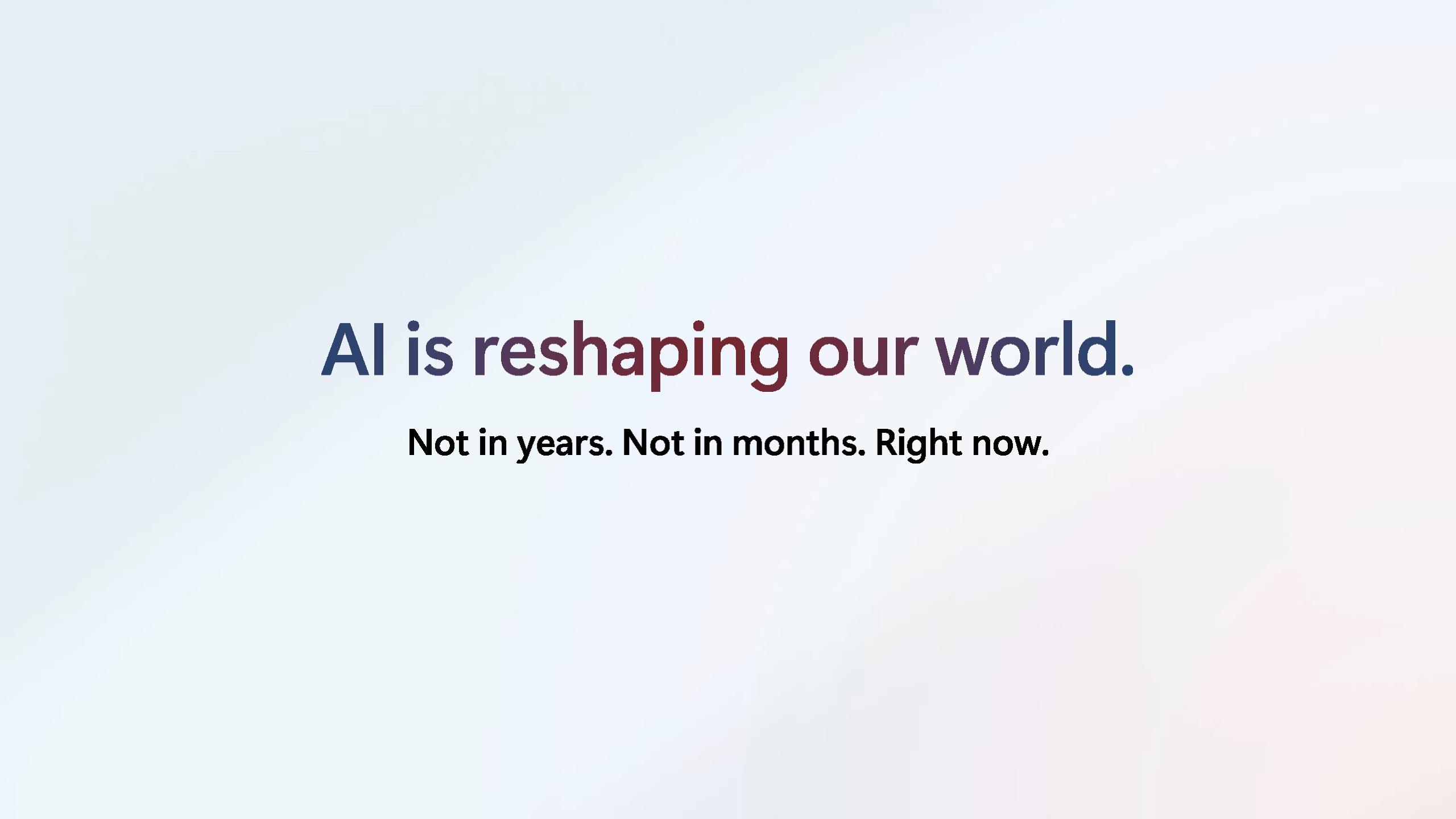
 <https://www.linkedin.com/in/linotadros/>



Distinguished executive leader and renowned technical expert in AI, Data, and Machine Learning. Leads cross-functional architectural teams to award-winning performance by developing strategic roadmaps and powering enterprise-wide projects. Serves as board member and advisor for multiple corporations, delivering strategic guidance on product line developments and business solutions. Industry influencer and mastermind of strategic programs and innovations, leading modernization efforts to alter the global IT landscape as a Microsoft MVP and Microsoft Regional Director for 23 years.

Lino has been in the industry for 33 years and has presented at technology conferences in 54 countries so far, with over a dozen books published and hundreds of publications.

# The future of AI



# AI is reshaping our world.

**Not in years. Not in months. Right now.**

# AI-driven business transformation



Enrich  
employee  
**experiences**



Reinvent  
customer  
**engagement**



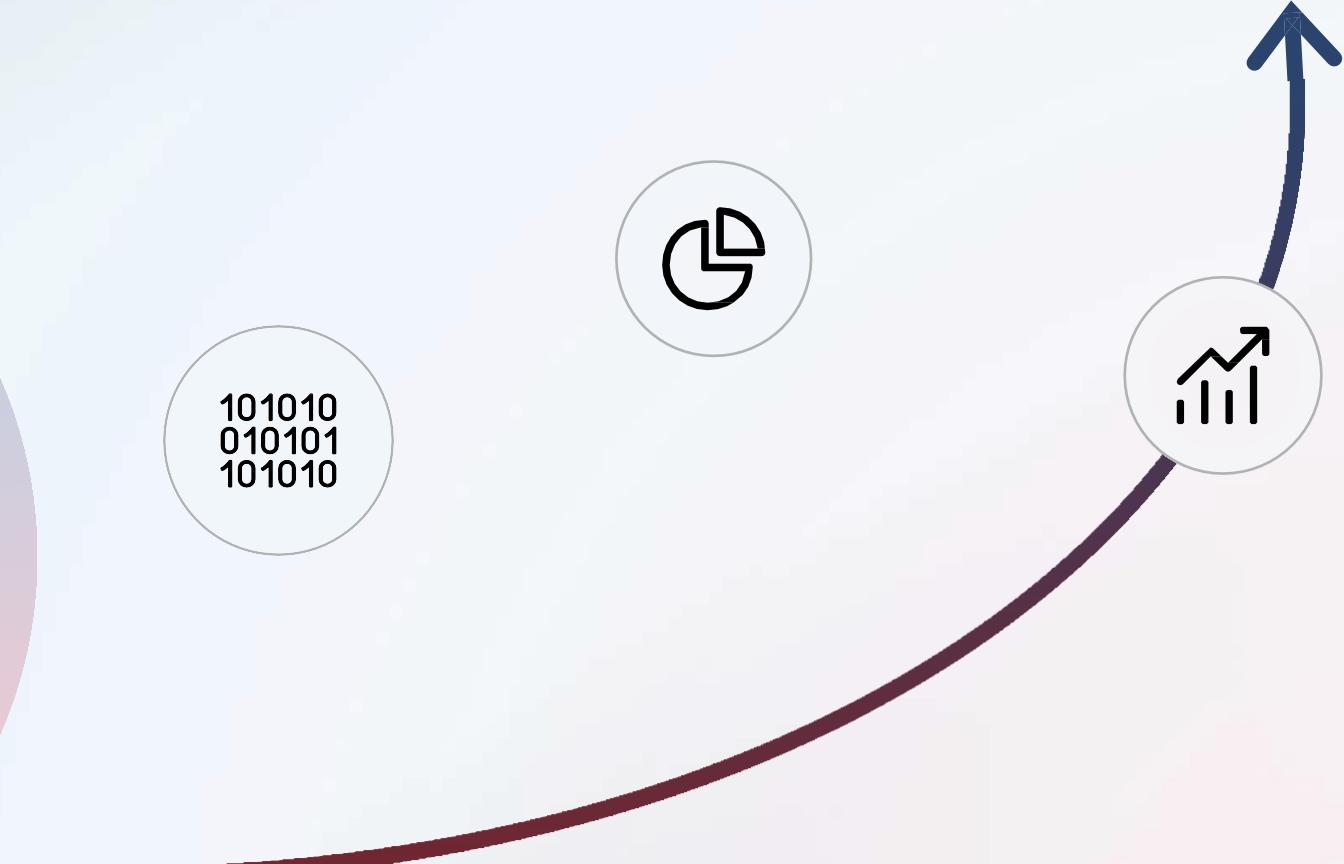
Reshape business  
**processes**



Bend the  
curve on  
**innovation**

# The business case for investing in AI

For every \$1 a company invests, the return on investment is \$3.7x.



# Generative AI trends

**81%**

of leaders expect agents to be integrated into their company's AI strategy in next 12-18 months<sup>1</sup>

**93%**

organizations are experimenting with multiple models<sup>2</sup>

**70%**

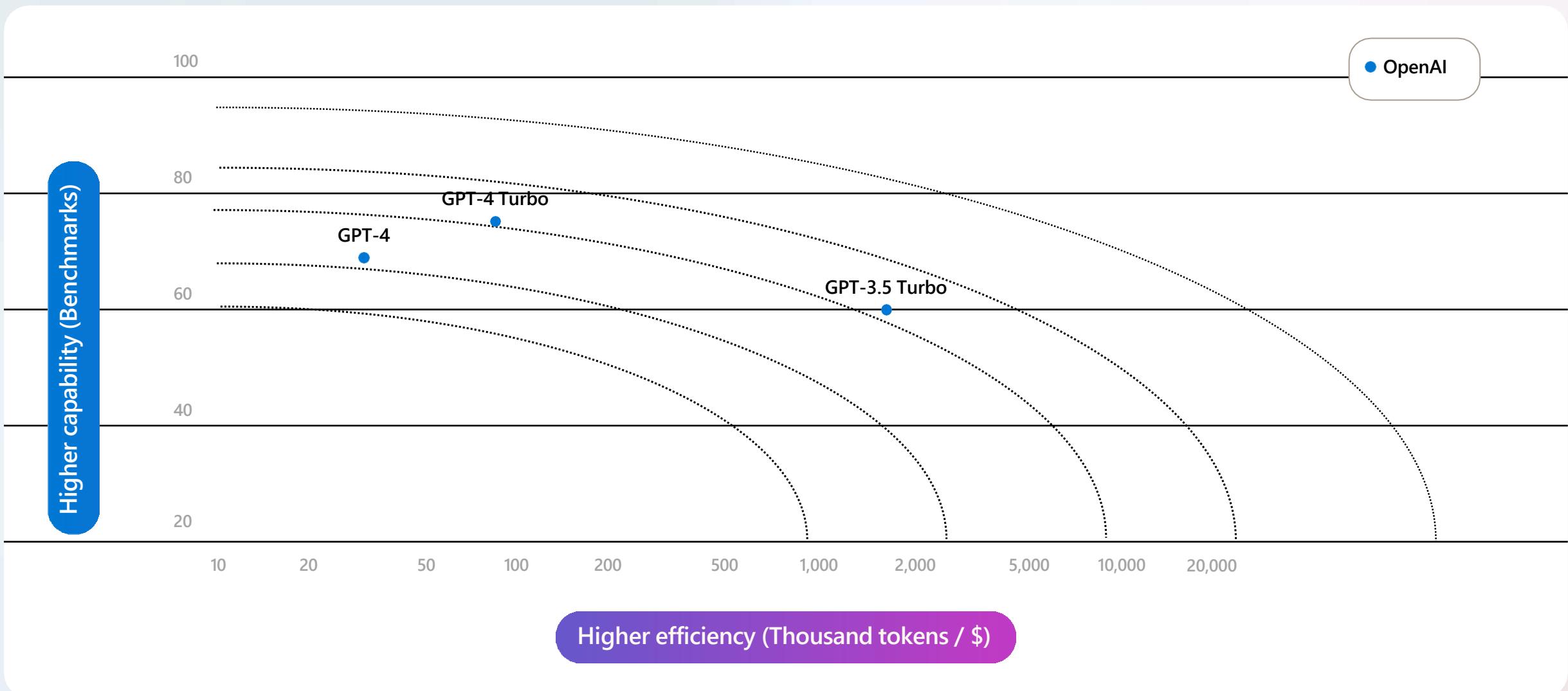
of generative AI experiments have not moved to production<sup>3</sup>

1. [Work Trend Index Annual Report](#)

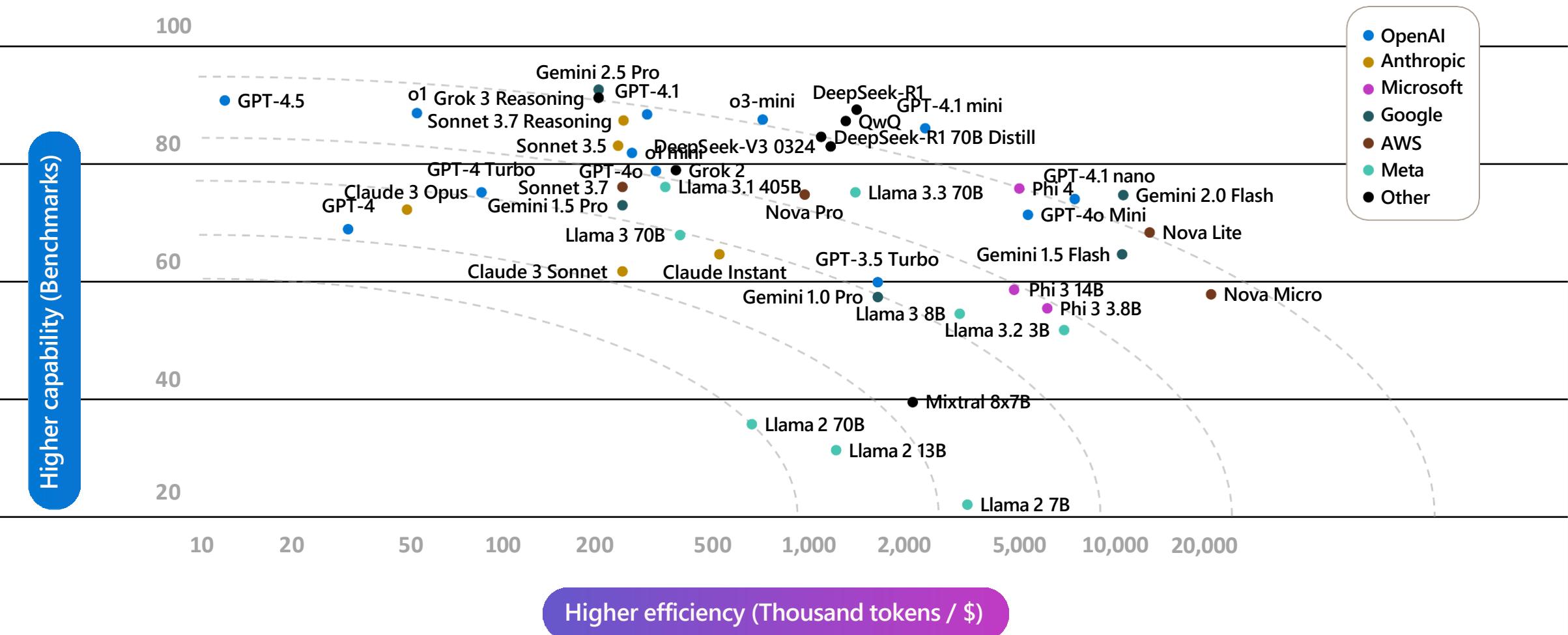
2. [16 Changes to the Way Enterprises Are Building and Buying Generative AI | Andreessen Horowitz](#)

3. [GenAI and the future enterprise | Deloitte Insights](#)

# From a few foundational models...

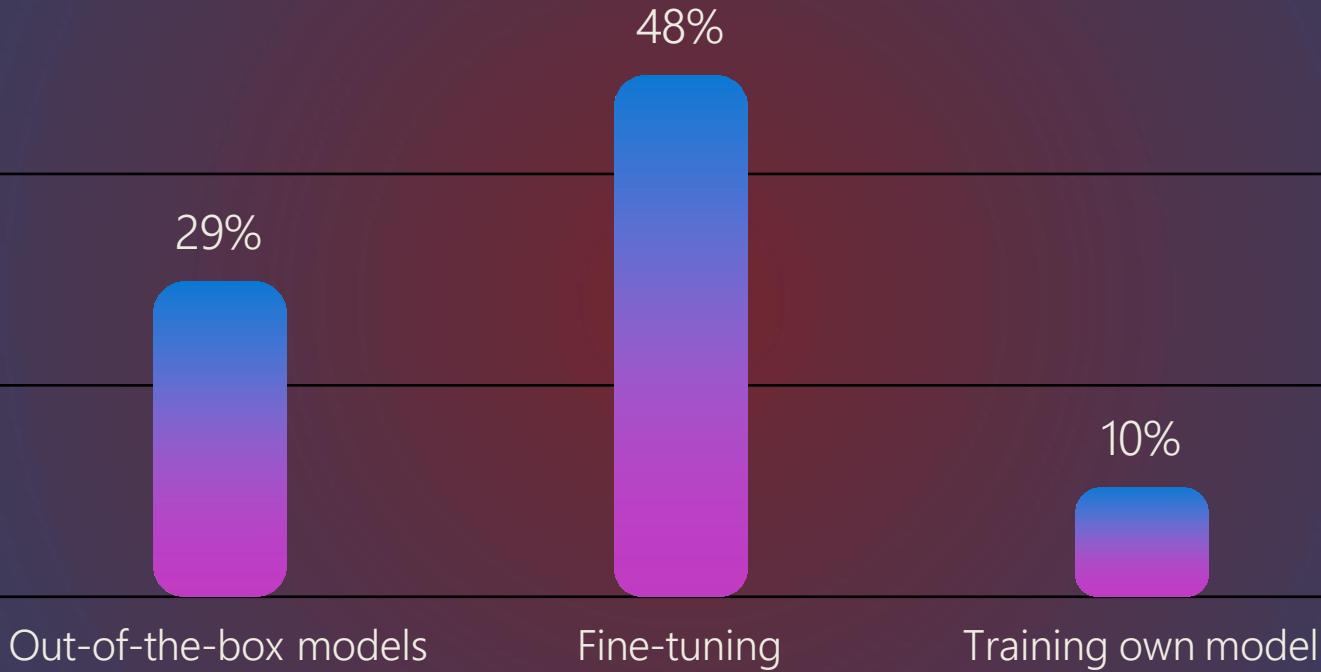


# ...to explosion of foundational models



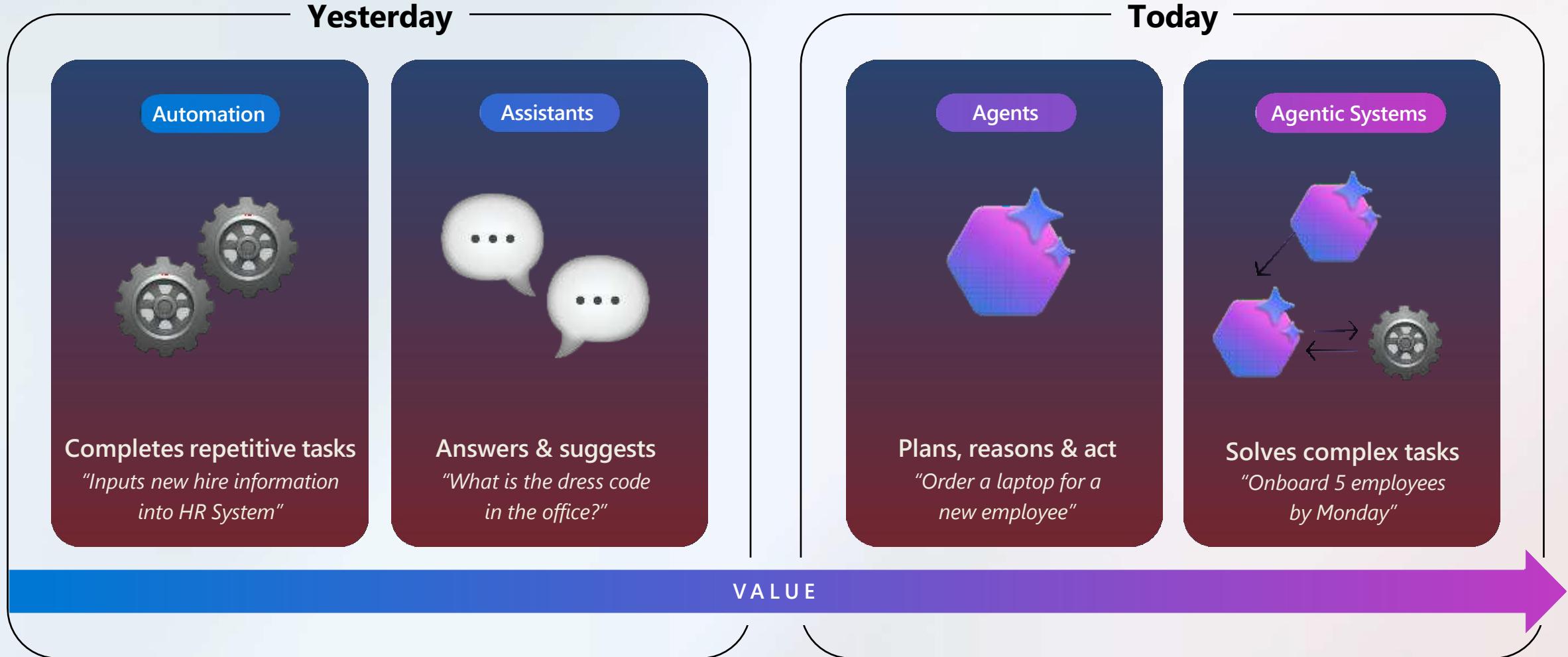
# Fine-tuning is becoming critical

Developer Approach to Models

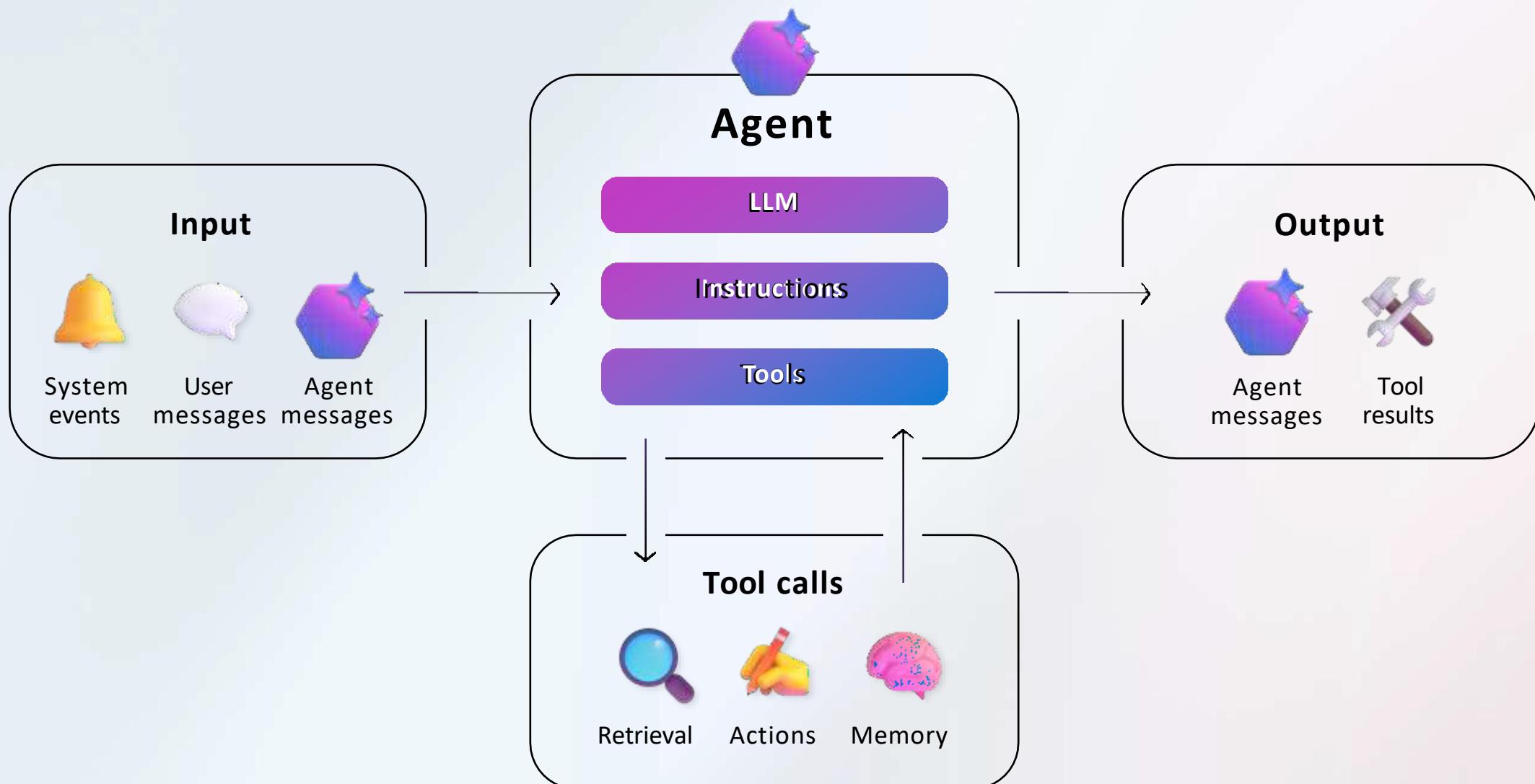


Source: Insight Partners 2024 Survey of enterprise tech leaders; 13% of the companies in the survey reported not building GenAI Apps

# Agentic AI is the new frontier



# What is an agent?



# Agents are already changing the way we...



Build and  
**manage AI**



Interact with  
**software**

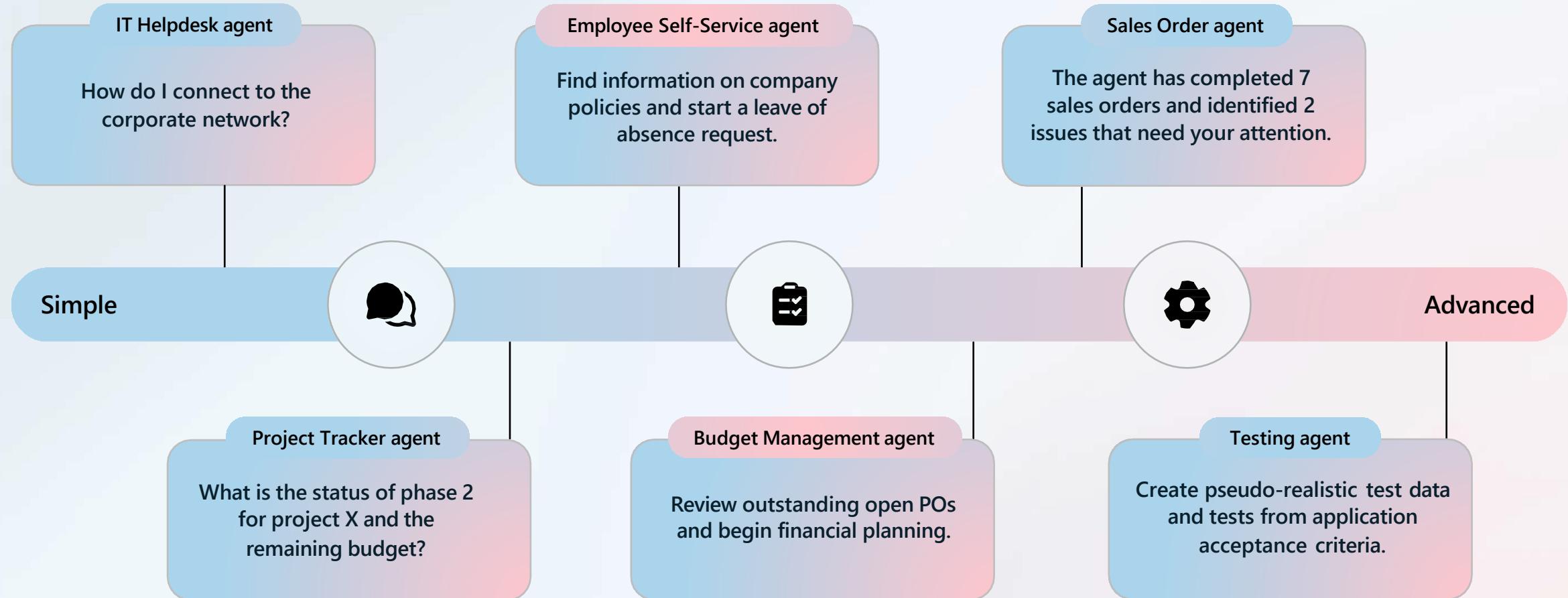


Gain value  
**from data**



Solve complex  
**challenges**

# ...with customized solutions...



# Who is using Generative AI?



# **...and the way we write code is changing**



**46%**

of code is now  
written by AI

Boosting developer  
productivity by

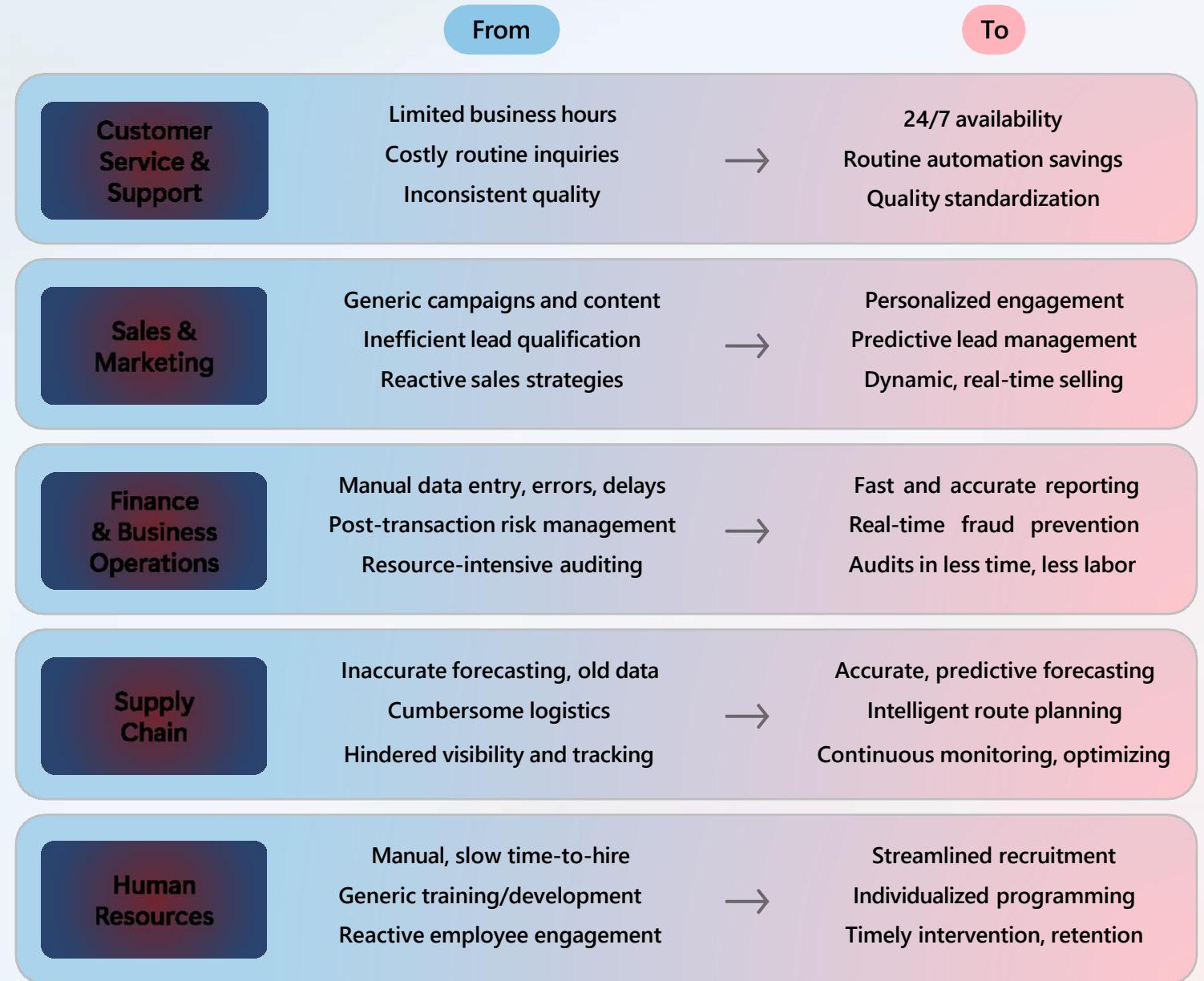
**55%**

AI tools

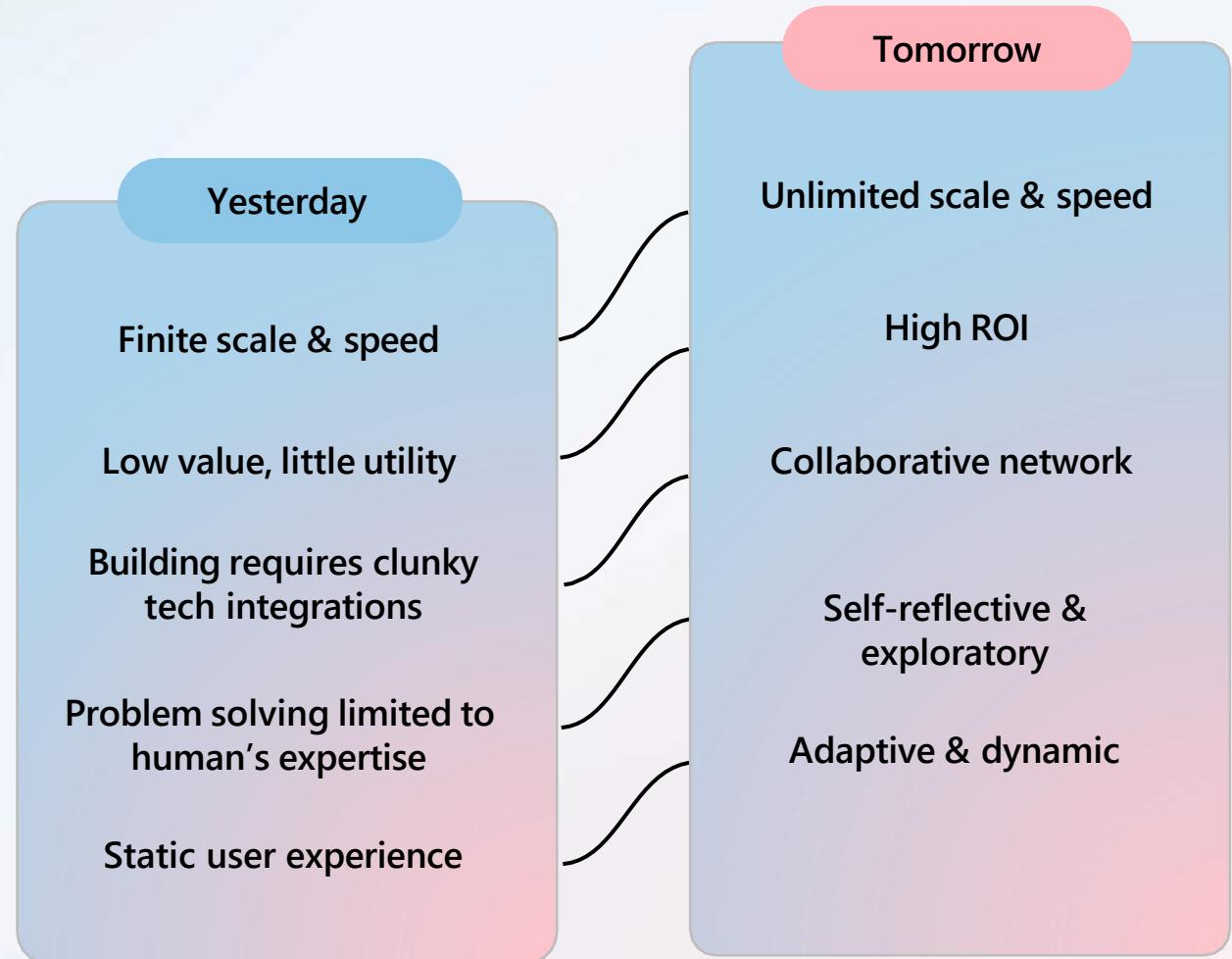
Vibe coding

Low code

# Every application achieves more with AI and agents



**...and they ask more  
of your applications**



# Agentic AI development challenges

70%  
or more generative AI  
experiments never  
make it to production<sup>2</sup>



Model selection  
& deployment



Complex  
workflows



Content  
safety



Observability  
and governance



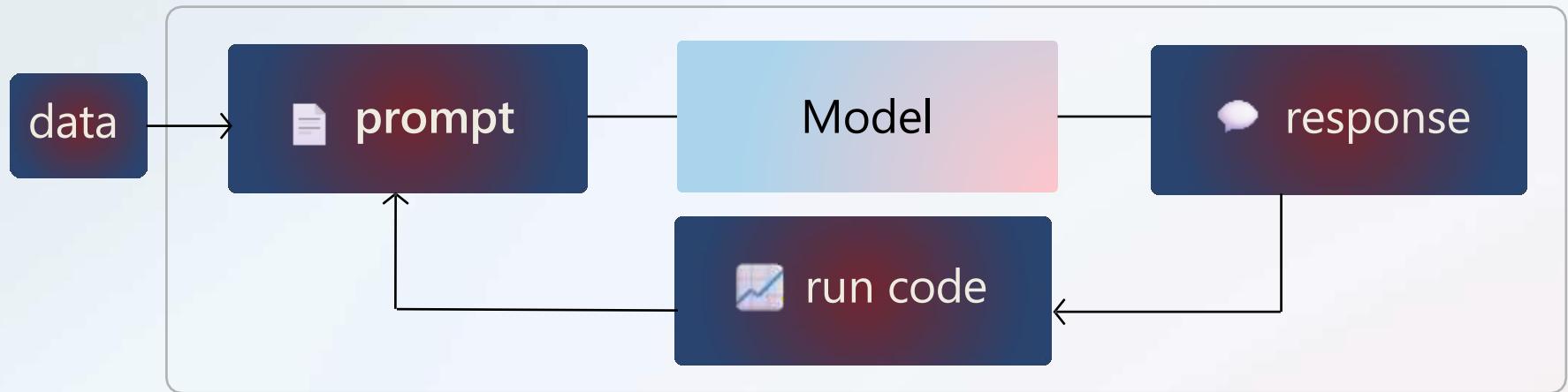
Tool  
sprawl

# De-mystifying AI App Development

## At a basic level

Providing prompt to a model

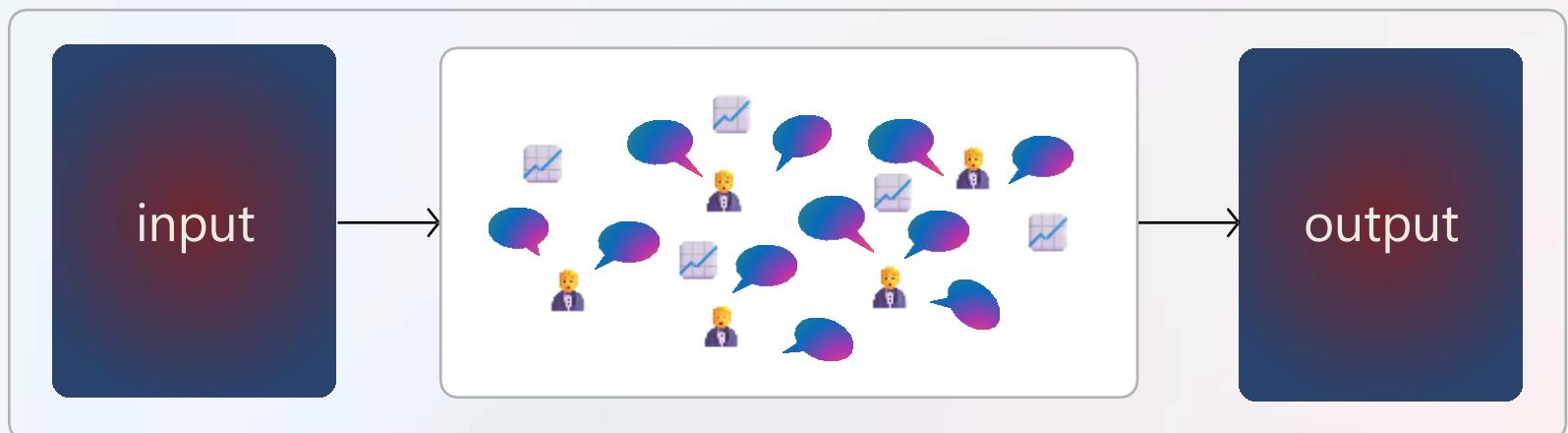
Acting on the response



## At scale

Holding many different conversations with models & agents

Extracting information into useful results



# Developer needs



Rapidly prototype,  
**build, and deploy** to  
share concepts with  
customers



Easy experimentation  
and evaluation to  
quickly test prompts,  
models, and datasets  
without waiting for  
provisioning



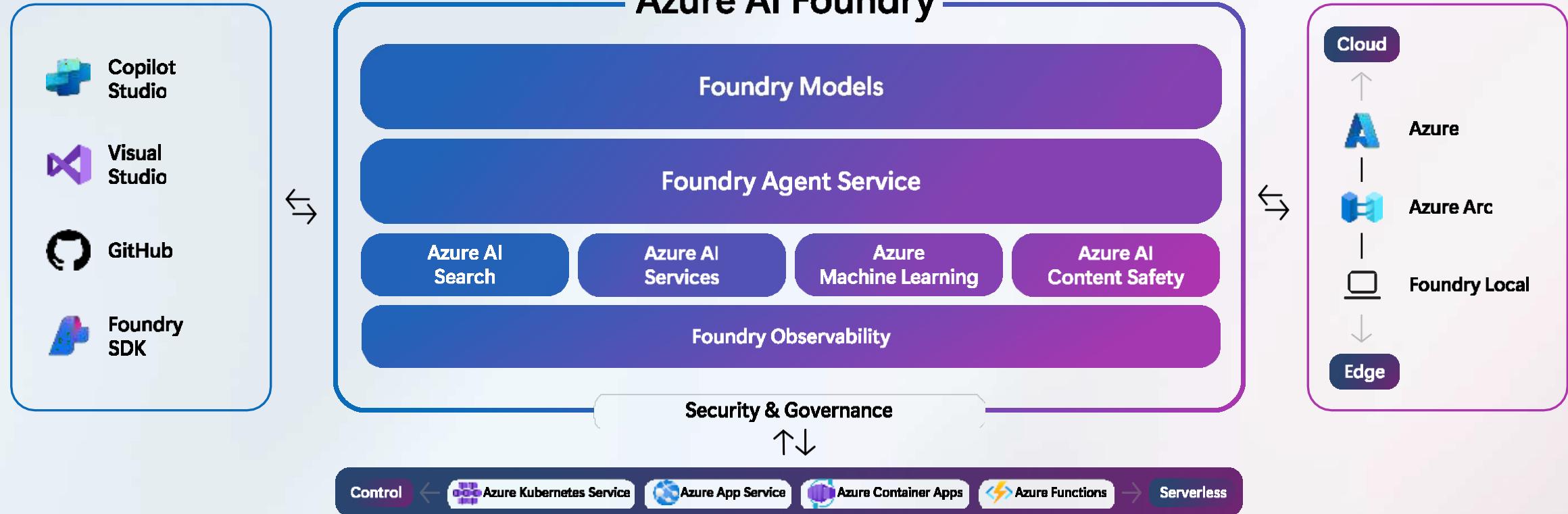
Secure, scalable,  
**production ready AI**  
platform to manage AI  
applications and agents

# The AI app & agent factory

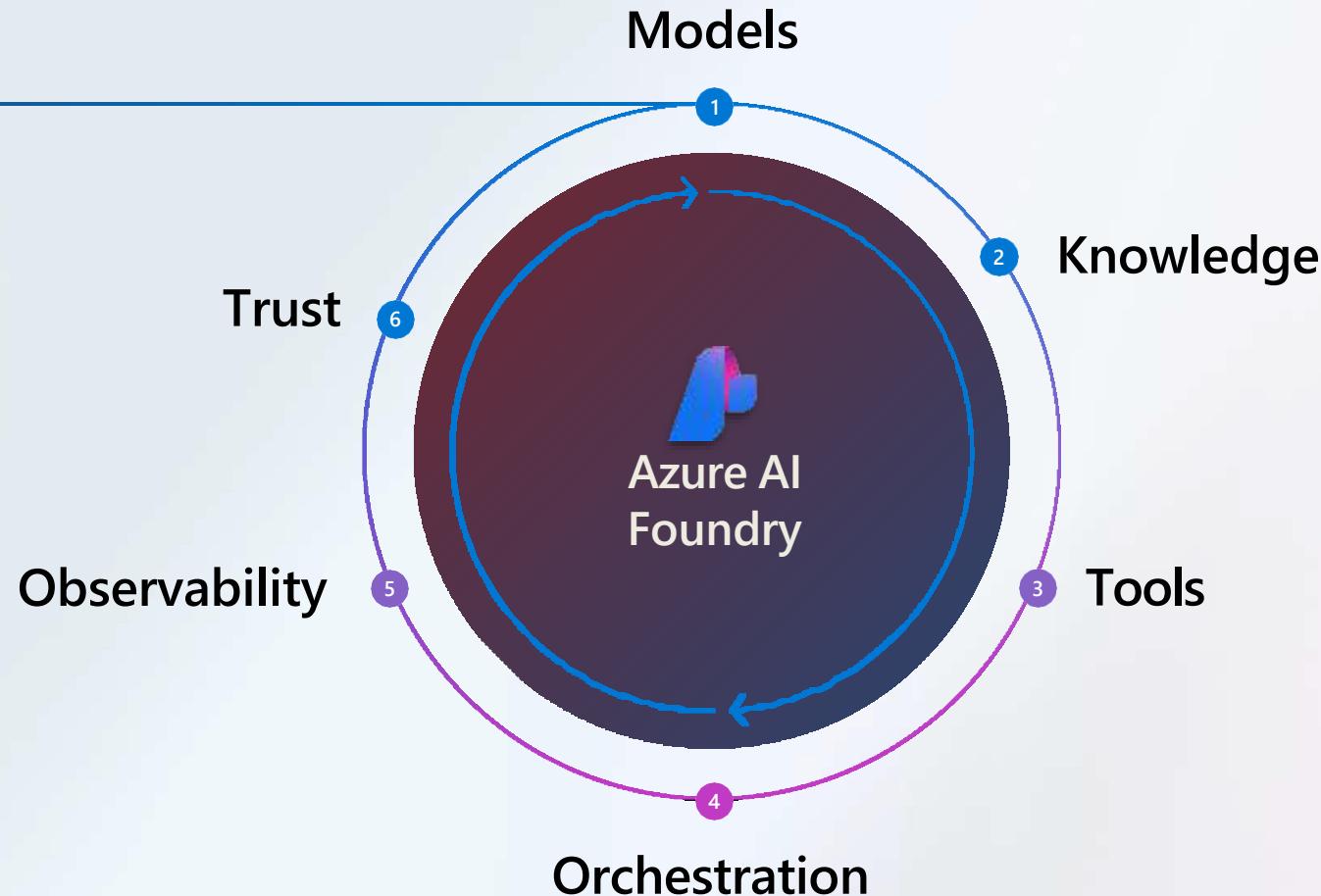
Design, customize, and manage AI  
apps and agents at scale



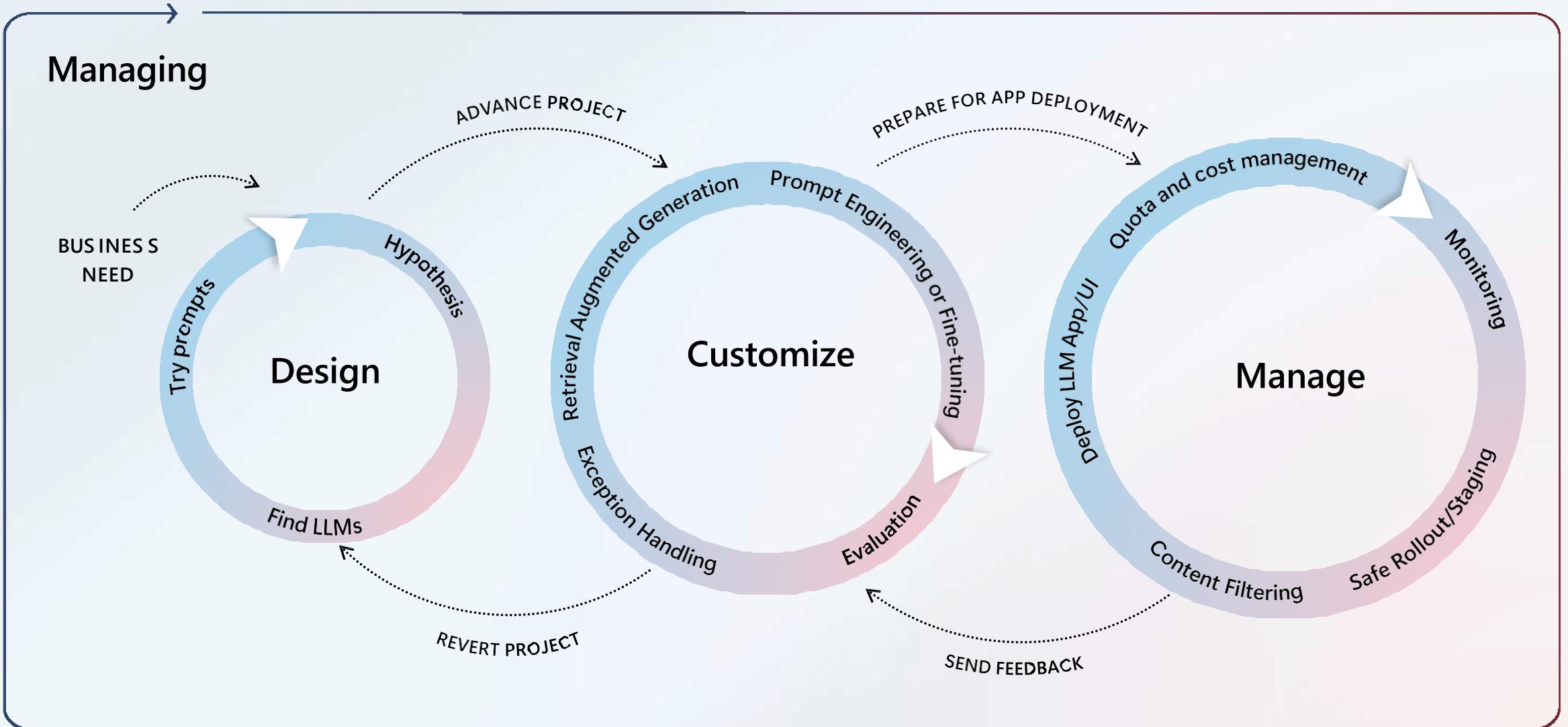
## Azure AI Foundry

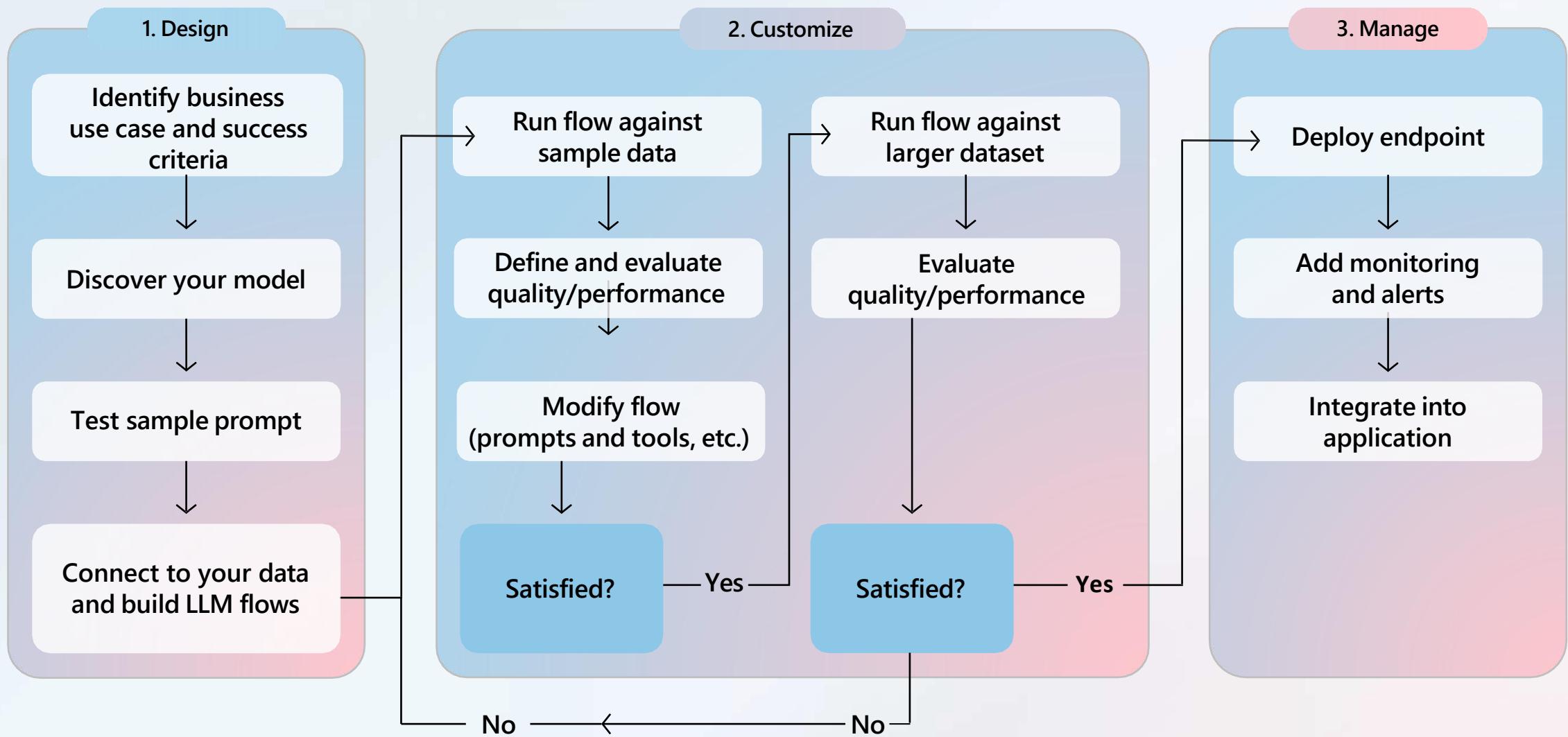


# The AI app and agent factory



# AI Development Lifecycle





# Microsoft runs on Azure AI Foundry

Microsoft  
Copilot

Microsoft 365  
Copilot

Microsoft Security  
Copilot

GitHub  
Copilot

DAX  
Copilot

# ROI with Azure AI Foundry



30% average  
productivity gain



95% reduction in  
configuration time



50% potential  
reduction in time to  
build AI applications



99.97% time  
reduction for  
model training



"We're anticipating more problems being reported and faster resolution time, and using [Azure AI Foundry] allows us to do that."

**Joe Bohman,**  
Executive Vice President,  
Siemens Digital Industries Software

## Business goal

Siemens wanted to bridge the gap between field and shop floor workers and operations and engineering teams to better drive innovation and efficiency, but more importantly, help companies rapidly address problems as they arise.

## Technology solution

Using Microsoft Teams and Azure AI Foundry, Siemens created an **AI-powered collaborative app** to help support their industry leading product lifecycle management (PLM) solution, Teamcenter, and connect people who find problems with those who can fix them.

## Impact

The Siemens Teamcenter for Microsoft Teams app promises to close feedback loops faster and help cross-functional teams resolve more issues faster—and together.

[Read full story here](#)

## Key products

Foundry Models



"We really appreciate the one-click deployment of the models in Azure AI [Foundry] and that it makes Azure AI offerings transparent and available to the user."

**Shah Muhammad,**  
Head of AI Innovation,  
Sweco

## Business goal

Sweco's engineers and architects manage numerous, complex tasks every day. To allow more time for more creative solutions in client projects, the company recognized the need for a supportive solution.

## Technology solution

The company chose Microsoft Azure AI Foundry to build its **custom copilot, SwecoGPT, a digital assistant** that automates document creation and analysis, delivers advanced search, and provides language translation.

## Impact

Though it was recently deployed, nearly half of Sweco's employees use SwecoGPT and report increased productivity, giving them more time to focus on creativity and helping customers.

[Watch the full story here](#)

## Key products

Foundry Models  
Azure AI Services  
Azure Machine Learning

# 70,000+ organizations using Azure AI Foundry today



Design with the best  
models for your use case

# AI templates in Azure AI Foundry



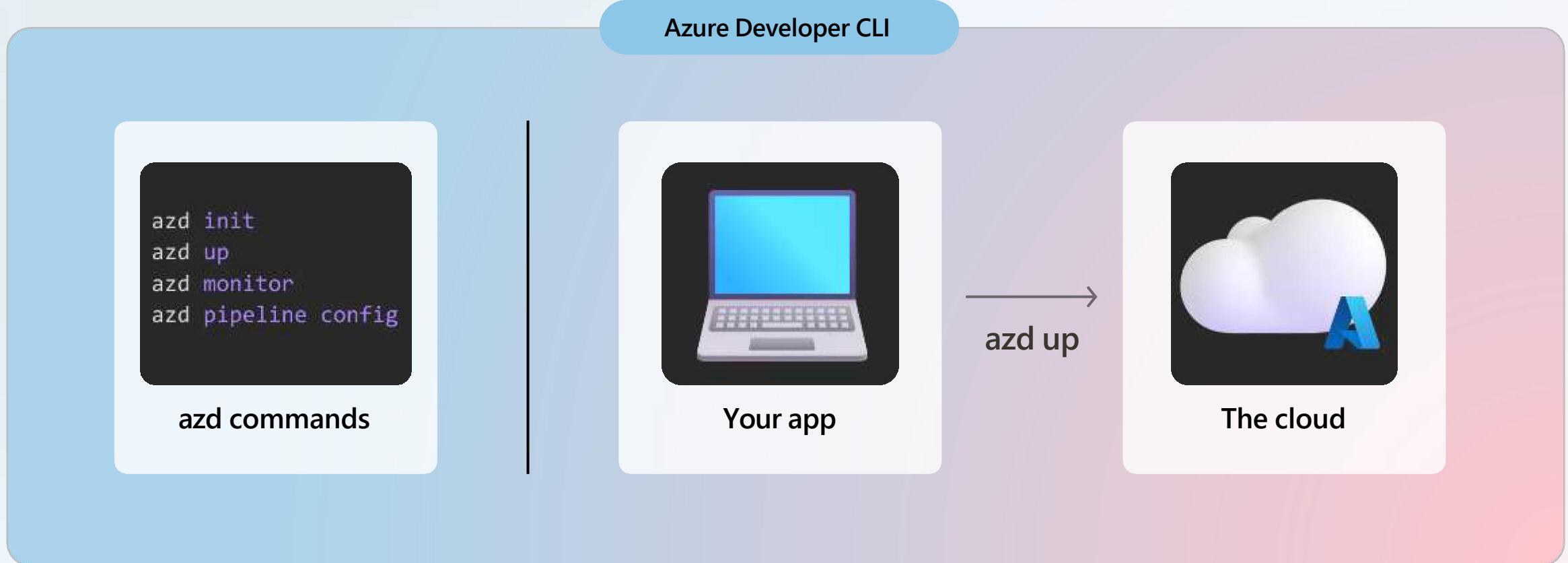
-  Unlock insights from conversational data
-  Get started with AI chat
-  Get started with AI agents
-  Create a conversational agent
-  Generate documents from your data
-  Improve client meetings with agents
-  Process multi-modal content
-  Modernize your code with agents
-  Automate workflows with multi-agents
-  Deploy your AI application in production

# AI templates in Azure AI Foundry

Get started quickly with customizable AI templates to support a wide variety of use cases.

Template	Popular Customer Use Cases
<b>Unlock insights from conversational data</b>	<ul style="list-style-type: none"><li>• Mining call center conversation insights</li><li>• Call center operations improvement</li><li>• Transcription</li></ul>
<b>Get started with AI chat</b>	<ul style="list-style-type: none"><li>• Employee chatbot for HR / benefits</li><li>• Professional Services assistance (Legal/Tax/Audit)</li><li>• Analytics and reporting</li><li>• Contact center agent assistance</li></ul>
<b>Get started with AI agents</b>	<ul style="list-style-type: none"><li>• Customer service</li><li>• Conversation summarization</li><li>• Translation</li></ul>
<b>Build your conversational agents</b>	<ul style="list-style-type: none"><li>• Client meeting prep</li><li>• Document generation</li><li>• AI-enabled product experiences</li><li>• Product discovery/shopping assistant</li></ul>
<b>Generate documents from your data</b>	<ul style="list-style-type: none"><li>• Clinician-patient visit documentation</li><li>• Claims/Contract/Invoice processing</li><li>• ID verification</li></ul>
<b>Improve client meetings with agents</b>	<ul style="list-style-type: none"><li>• Code conversion</li><li>• Code validation</li><li>• Documentation</li></ul>
<b>Multi-modal content processing</b>	<ul style="list-style-type: none"><li>• Employee On-Boarding and self-service</li><li>• Travel Booking and Expense Management</li><li>• Supply Chain Planning</li></ul>
<b>Modernize your code with agents</b>	
<b>Multi-agent workflow automation</b>	
<b>Deploy your AI application in production</b>	Taking proof of concept applications to production

# Local dev environment to the cloud in a single step



**Difficulty  
selecting models  
and services**

**New models  
launching  
and evolving**

**Need to run AI  
solutions  
everywhere from  
cloud to edge**

**AI innovation requires model flexibility**



## Azure AI Foundry Models

# Design with the best models for your use case

- Discover models for your use case
- Simplify model selection and upgrades
- Run anywhere - cloud, local devices and on-premises

The screenshot shows the Azure AI Foundry Model catalog interface. At the top, there's a navigation bar with links for 'All hubs + projects' and 'Project controls'. Below the navigation is a search bar and a main title 'Find the right model to build your custom AI solution'. The interface includes sections for 'Announcements' (with cards for 'New Phi reasoning models', 'Introducing gpt-image-1', 'MAI-DS-RT reasoning model', and 'Unlock Enterprise Agent workflows with o3 and o4-mini'), 'Model leaderboards' (showing top models for Quality, Cost, and Throughput), and a large grid of model cards. Each card displays the model name, icon, and description. A footer at the bottom reads 'Management center' and 'Image may not reflect actual user interface.'

# Azure AI Foundry Models

Offering 11,000+ frontier and open models

Built-in safety

**Deployment Types:**  
Standard  
Provisioned Throughput  
Batch (50% off)

**Deployment Locations:**  
Global  
Data Zone (US and EU)  
28 Regions

**Sold by:**  
Microsoft

## Hosted and sold by Microsoft



OpenAI  
Model Family  
*(available day 1)*



DeepSeek  
latest



Mistral AI  
latest  
*(coming soon)*



Meta Llama  
latest  
*(coming soon)*



xAI  
Grok  
*(coming soon)*



Black Forest  
Labs  
*(coming soon)*



Microsoft Phi  
Family

## Other Foundry Models



Hugging Face



NVIDIA  
NIMS



Cohere

## Specialized & Industry models



Unified access

Scalable deployment

Enterprise-ready

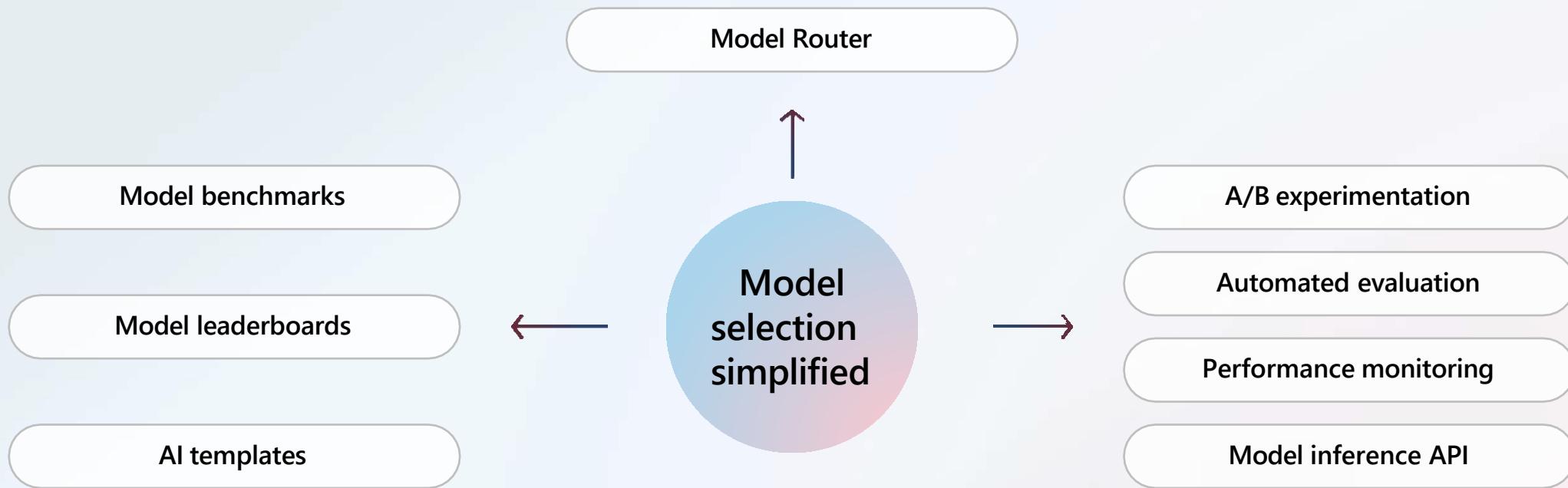
**Deployment Type:**  
Standard  
Managed Compute

**Deployment Locations:**  
Global  
7 Regions

**Sold by:**  
Partner

# Model selection simplified with Azure AI Foundry

Use the best model for your use case

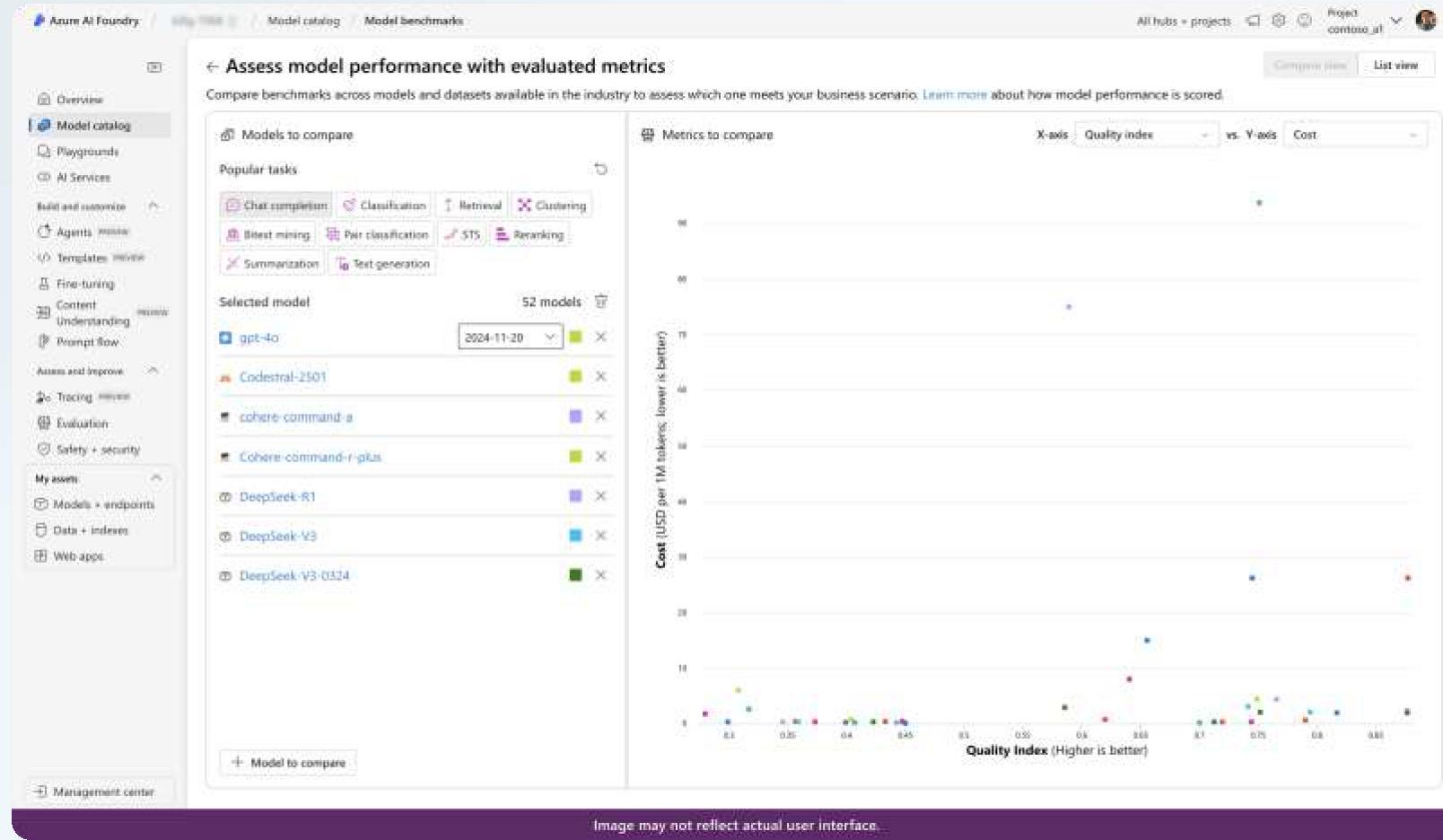


Select the best model

Continue using the best model

# Benchmarking in Foundry Models

Simplify model selection. Review and compare the performance of AI models.



# Leaderboard in Foundry Models

Identify top-performing AI models across different categories and tasks.

The screenshot shows the Azure AI Foundry interface for managing AI models. The left sidebar includes links for Home, Model catalog, Model leaderboards (which is selected), Playgrounds, AI Services, Build and customize (Code, Fine-tuning, Prompt flow), Assess and improve (Tracing, Evaluation), Safety + security, My assets (Models + endpoints, Data + indexes, Web Apps), and Management center. The main content area is titled "Model leaderboards" and displays the "Quality leaderboard". It lists the top 5 models based on Quality index:

Ranking	Model	Quality index
No. 1	AzureOpenAI-o1-preview	97
No. 2	AzureOpenAI-o1-mini	96
No. 3	Phi-4	92
No. 4	Minstral-3B	88
No. 5	Cohere-embed-v3-multilingual	85

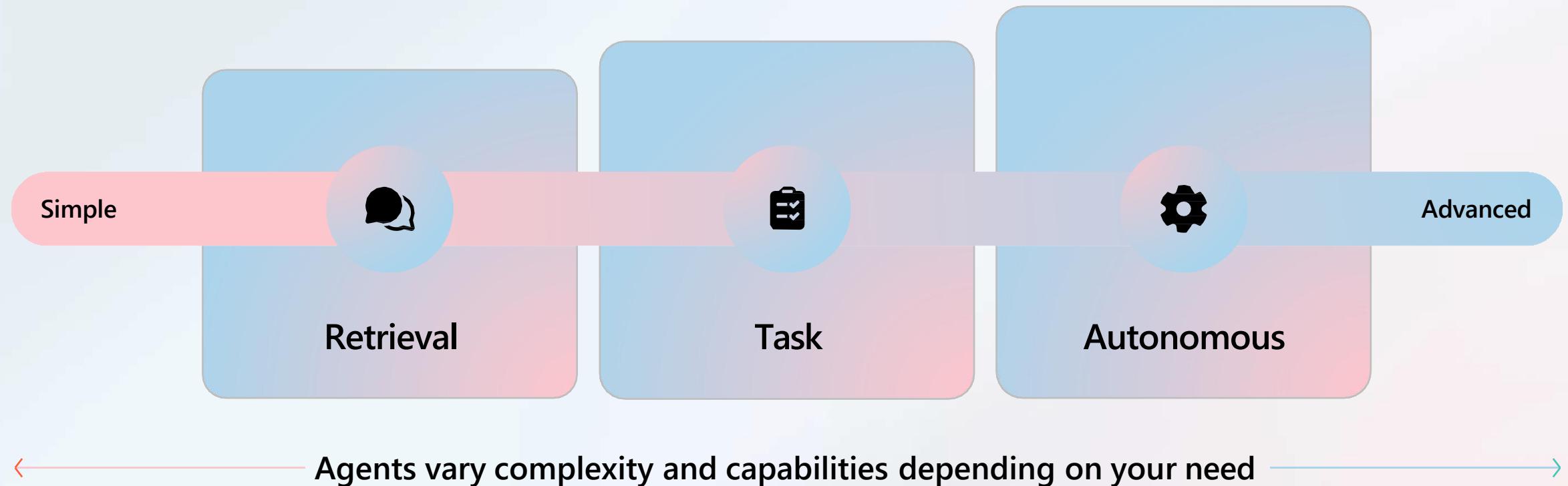
Below the table, there's a section titled "Highlights" showing three bar charts: Quality (Index: Higher is better), Speed (Output Tokens per Second; Higher is better), and Cost (USD per 1M Tokens; Lower is better). A callout box for the top model, "AzureOpenAI-o1-preview", provides more details: "Top of leaderboard for 3 weeks", "Chat completion, Question answering, Summarization", "Try in playground", "Deploy", and "Model details". It also notes "Deployed last month" and shows a line graph of usage with the value "6,273".

Image may not reflect actual user interface.

# AI agents

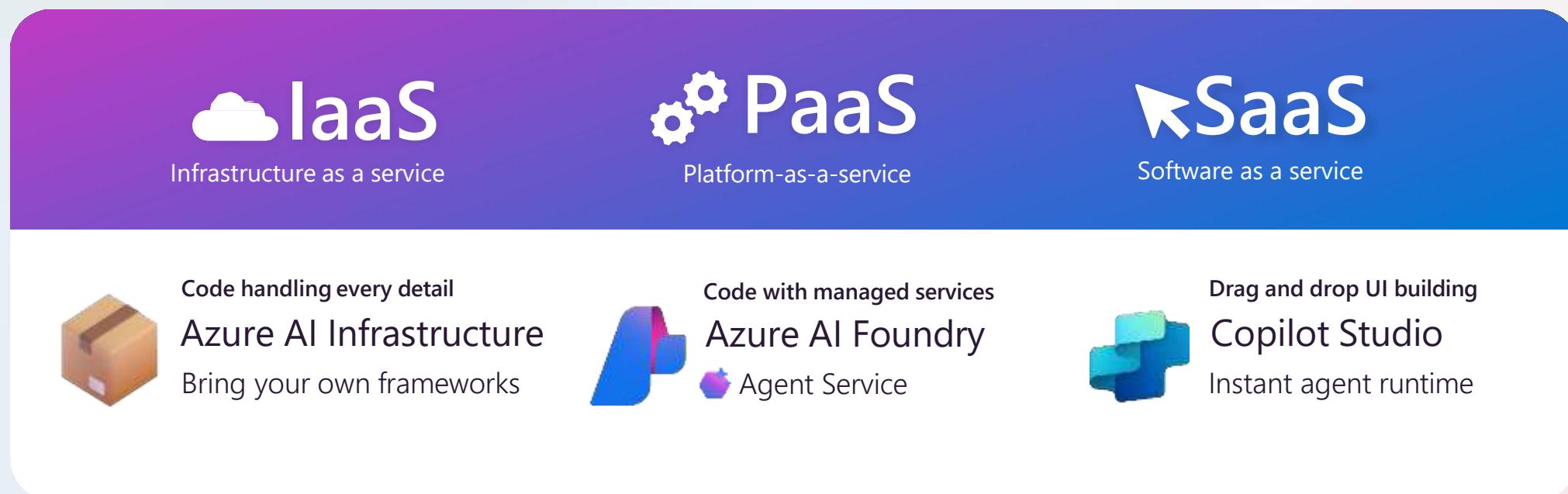
# What are agents?

Agents are programs that use AI to automate and execute business processes, working alongside or on behalf of a person, team or organization



# Freedom to create agents your way

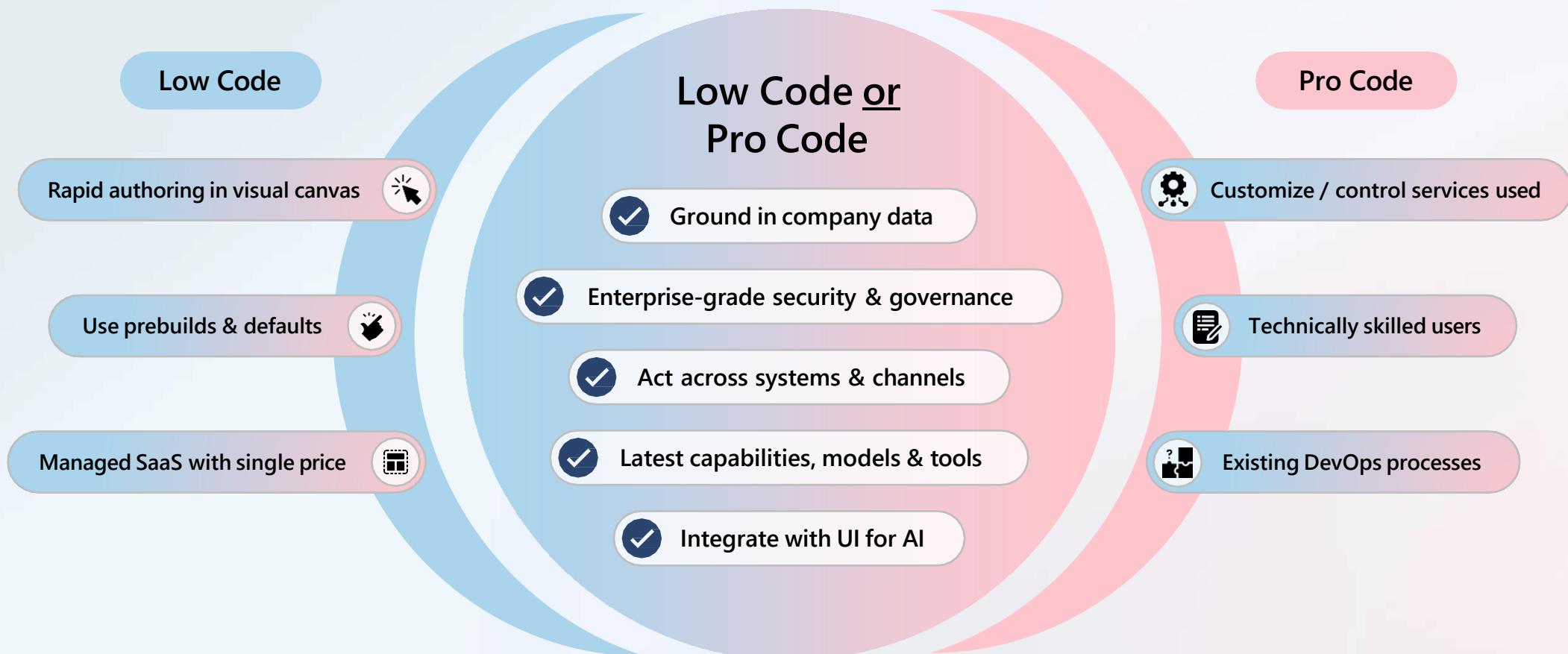
Platform Integrations, ease-of-use, and development speed

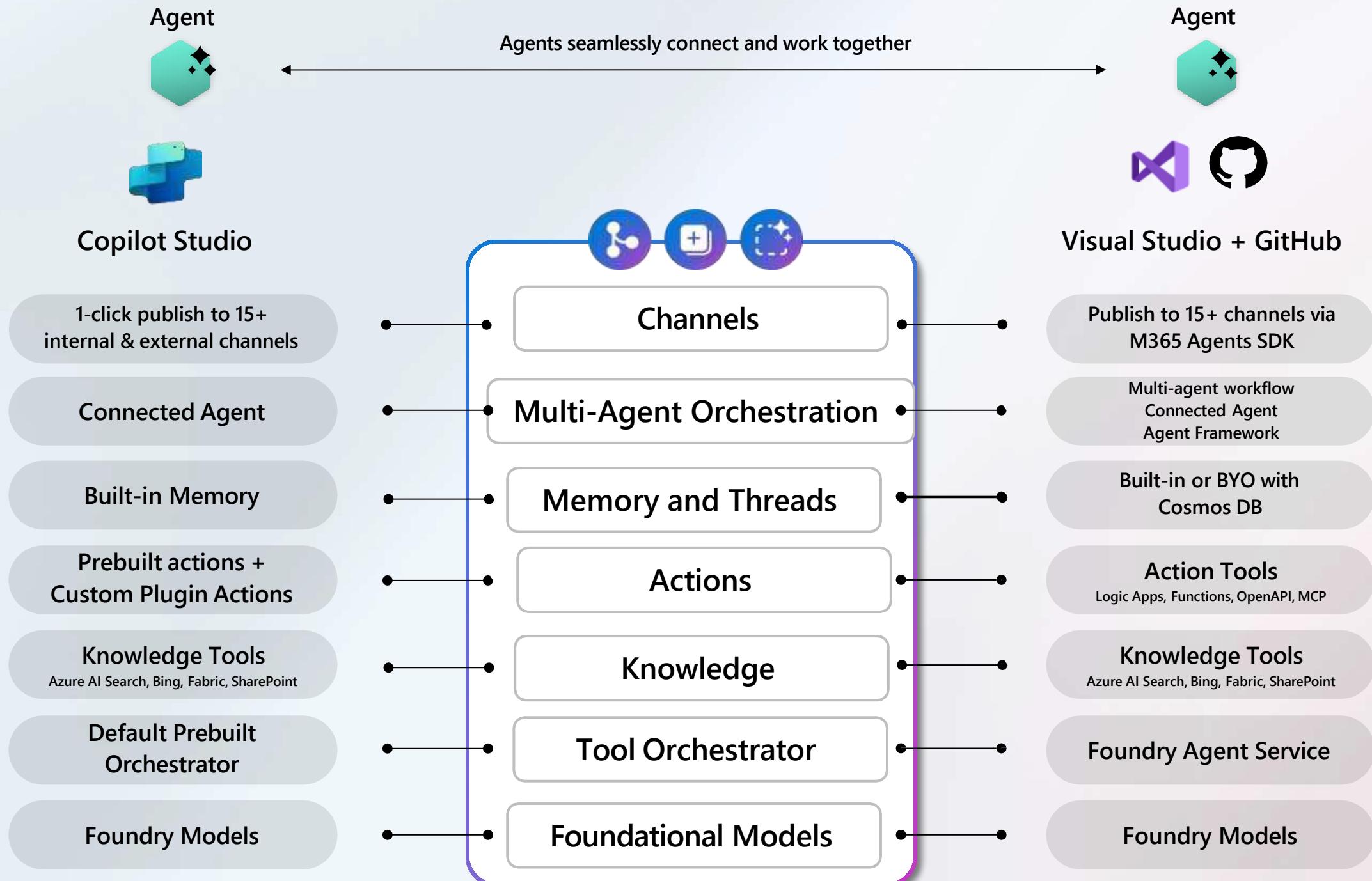


Control, visibility, and customization

# Choose the tool that best fits your needs

Most agents can be built using either of these differentiated tools  
– powered by Azure AI Foundry





# Securely developing AI agents requires 4 primary considerations



## Knowledge

providing agents with the right context



## Actions

Giving them the tools to complete their tasks



## Security

Provide only access to what they should



## Evaluation & Control

Check they complete their tasks correctly



# Azure AI Foundry Agent Service

Securely customize, orchestrate, and deploy AI agents

## Models

Model choice and flexibility  
with Foundry Models



Azure OpenAI

o1, o3-mini, GPT-5, 4.1, 4o, etc



### Azure Direct Models



Llama 3.1-405B-Instruct



Mistral Large



Cohere-Command-R-Plus

## Tools

Richest set of enterprise  
connectivity

### Knowledge



### Actions



Logic Apps\* Azure functions OpenAPI

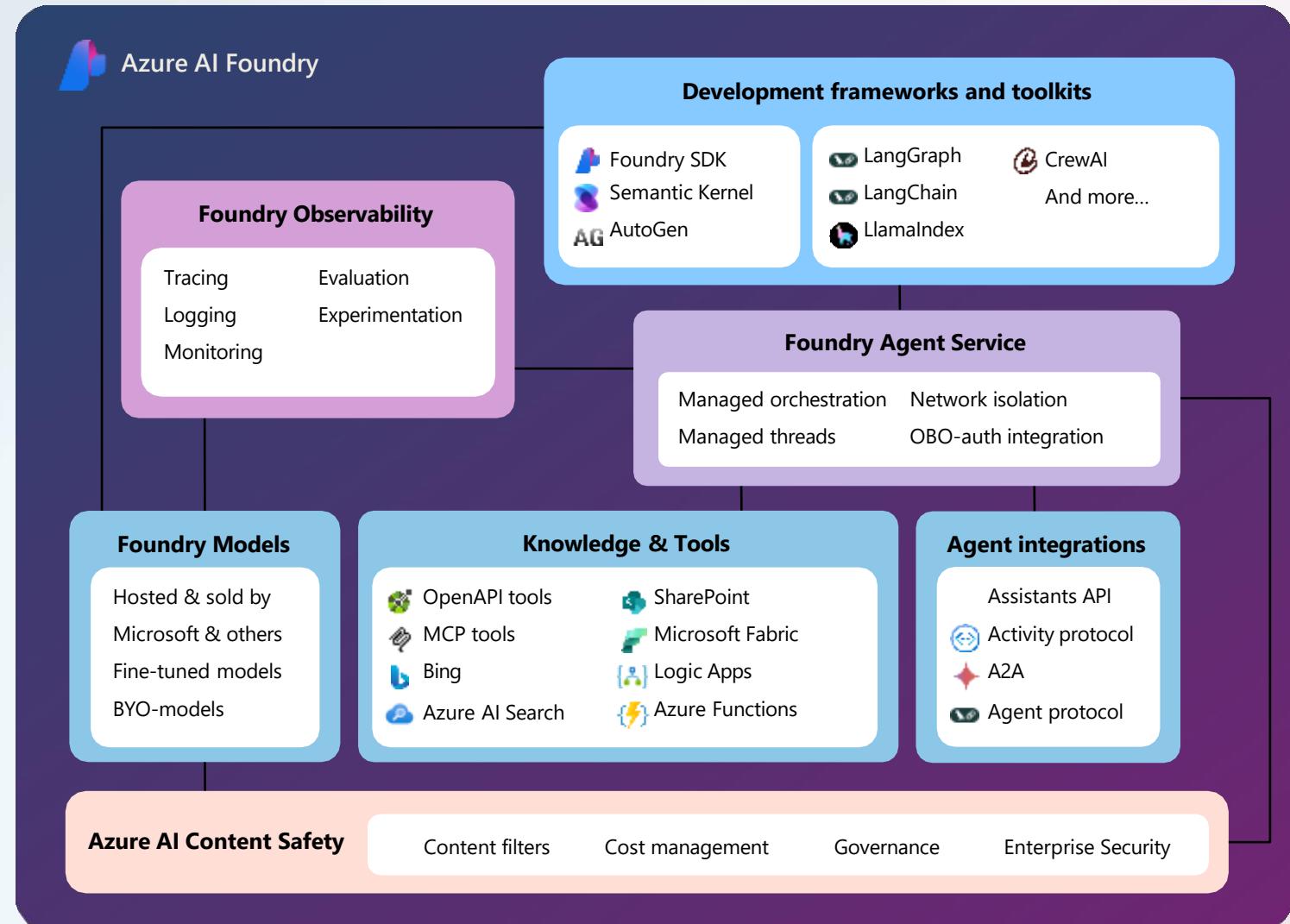
## Trust

Customer control over data,  
networking, and security

- BYO-file storage
- BYO-search index
- BYO-virtual network
- BYO-thread storage
- OBO authentication
- Content filtering

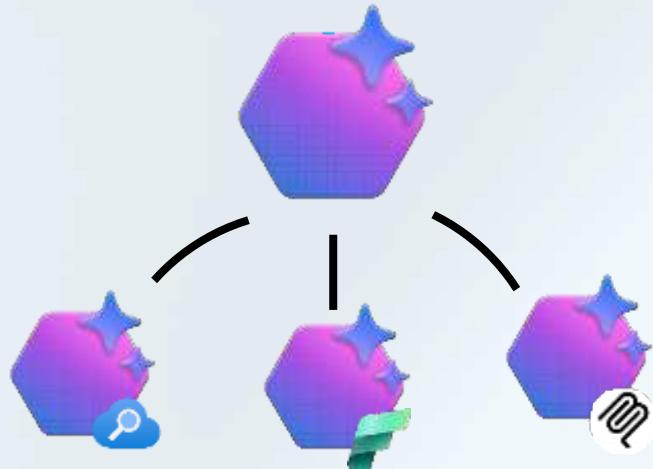
# Agents in Azure AI Foundry

*Combine the best models, services, and tools in Azure AI Foundry into reusable, testable agentic components.*



# Multi-agent orchestration in Foundry Agent Service

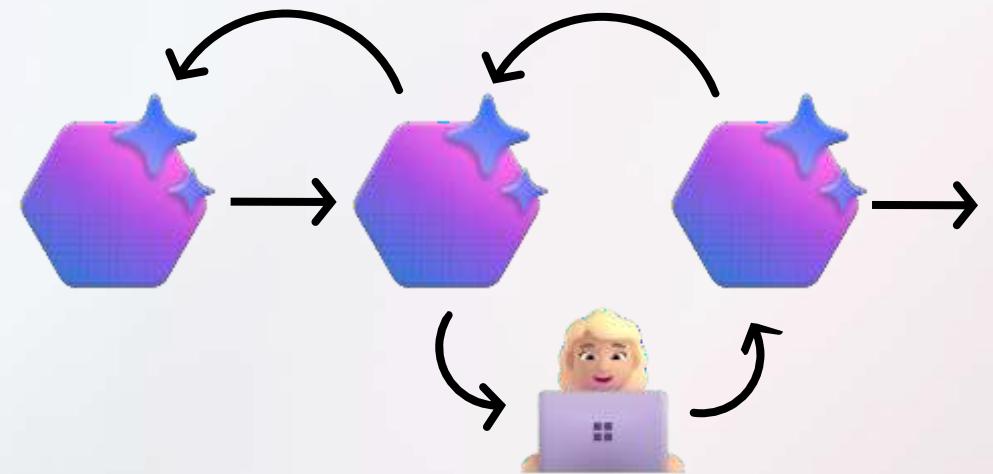
## Connected agents



Give one agent the abilities of another.

## Multi-agent workflows

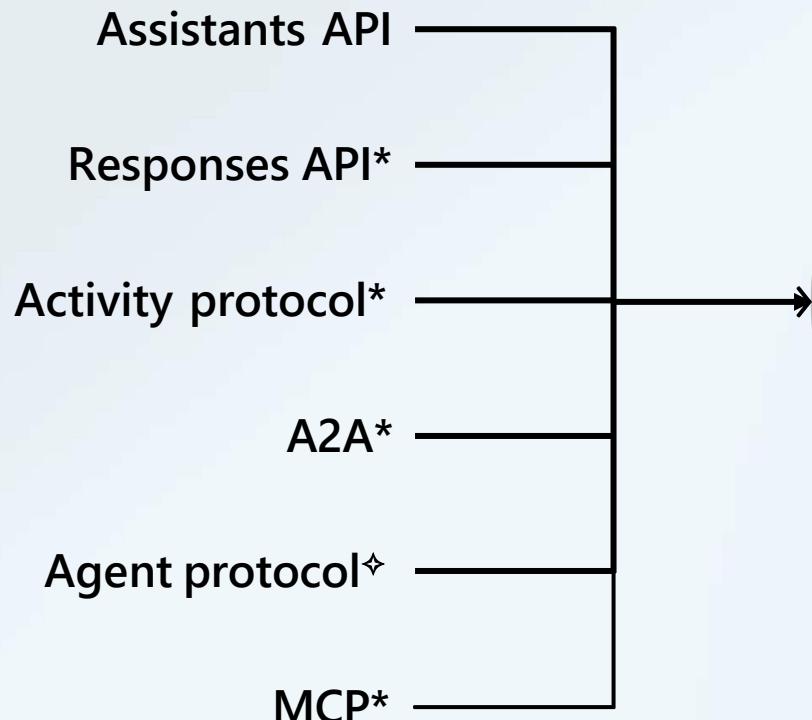
Powered by Semantic Kernel



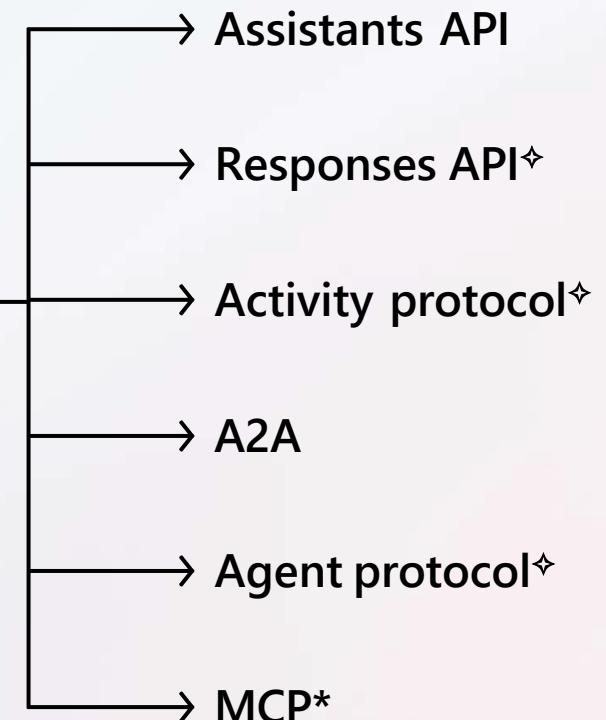
Declaratively orchestrate multiple agents.

# Foundry Agent Service will work with the top agent protocols

## Import external agents



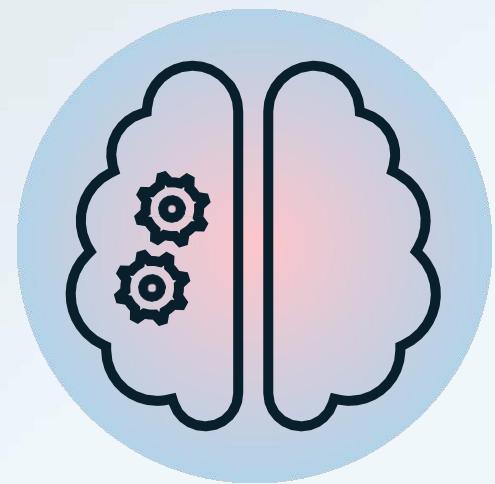
## Expose your agent externally



# Fine-tuning

# What is fine-tuning?

Customizing a pre-trained foundation model with additional training on curated data of a specific task or domain for enhanced performance, new skills, or improved accuracy.



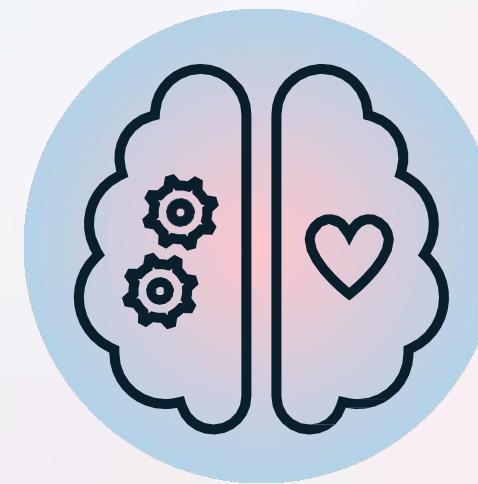
Generic LLM/SLM

+



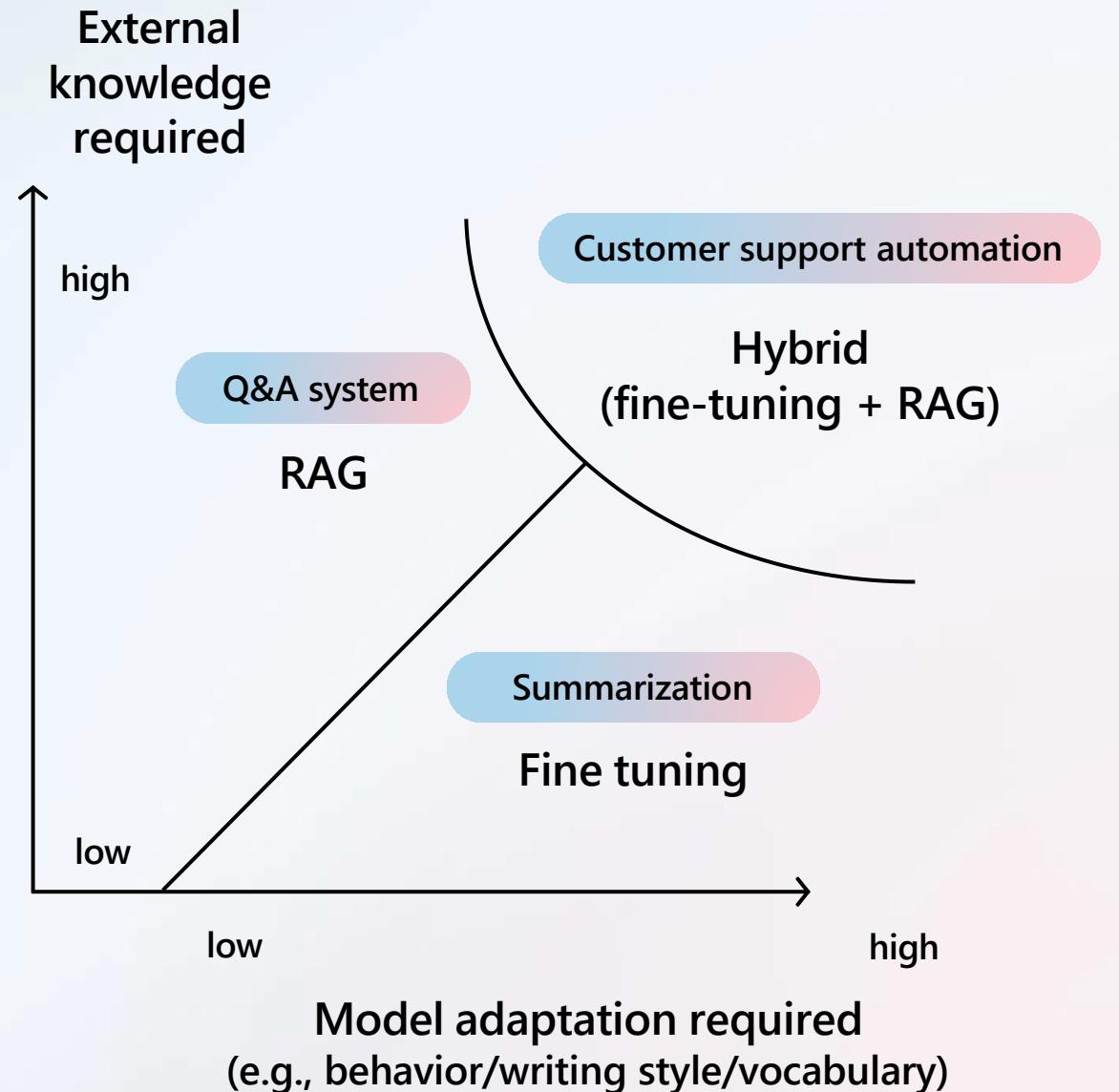
Curated Data Set

=



Fine tuned LLM/SLM

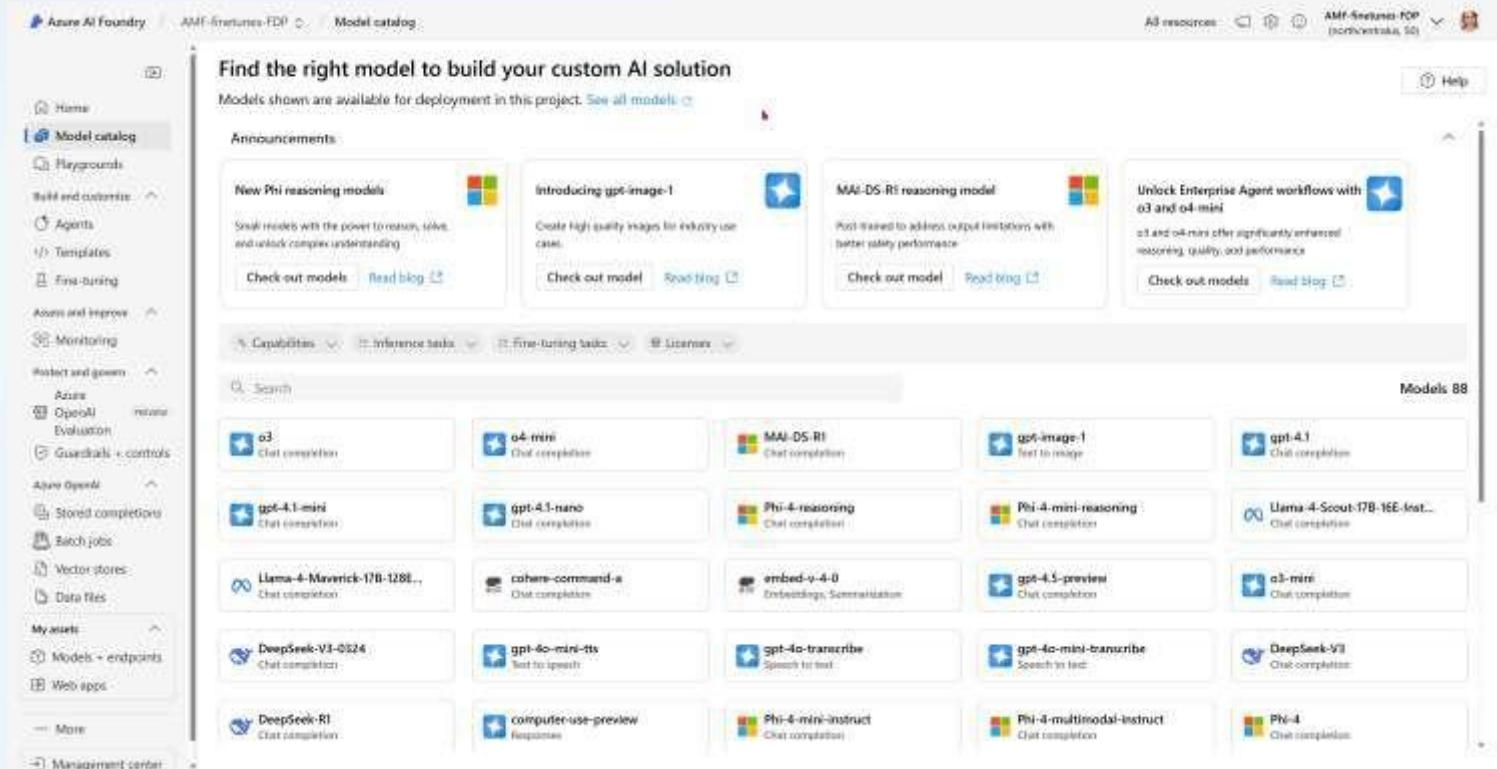
# When to use RAG vs. fine-tuning?



# Fine-tuning in Azure AI Foundry

Adapt Foundry Models for a specific task or on specific datasets to support business needs.

- Enhance performance and accuracy for specific use cases
- Reduce consumption costs and improve AI efficiency
- Lower latency
- Reduce bias and avoid hallucinations



Manage enterprise-ready  
AI with confidence

**Understand  
performance trade-  
offs for agentic  
applications**

**New attacks  
emerging and bad  
actors also using  
generative AI**

**Privacy, security,  
and compliance  
complexity**

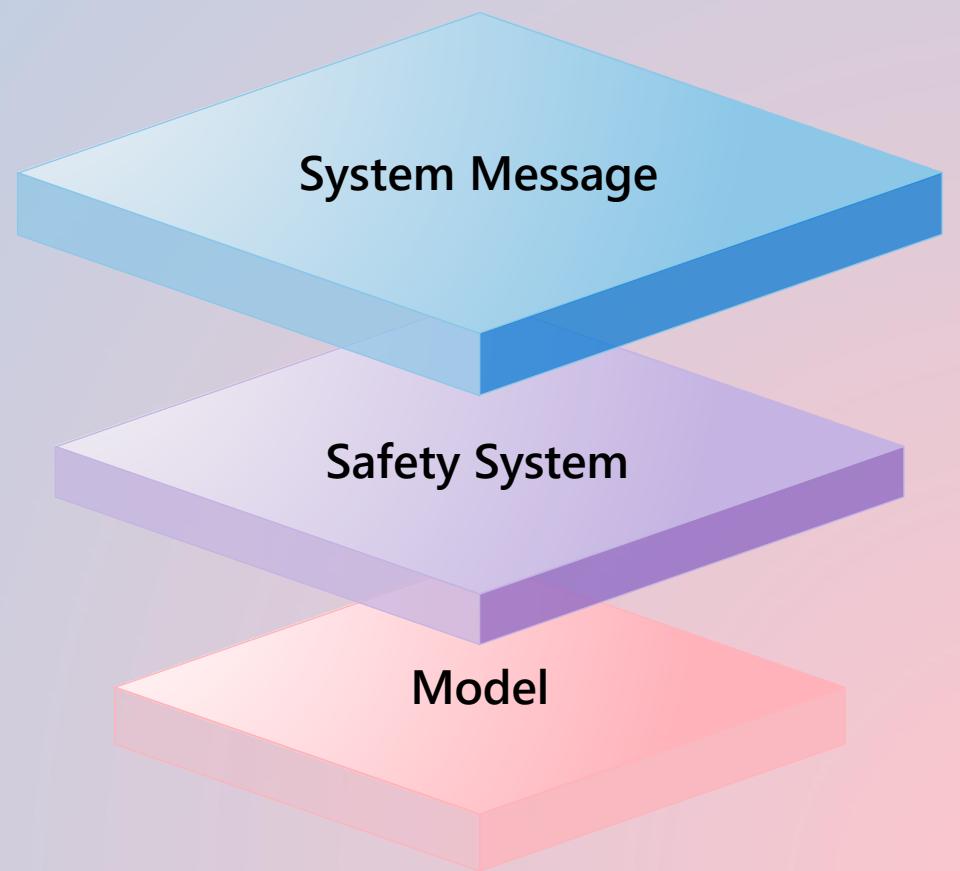
**Generative AI introduces new risk and needs to be managed with care**

# Mitigation layers in Azure AI Foundry

The **system message layer** provides hidden instructions to your model with every user prompt, so you can guide the model's behavior and data retrieval to generate higher quality responses by default.

The **safety system layer** acts as a shield or firewall around your foundation model, by detecting and filtering/blocking risky prompts before they reach your model and risky responses before they reach your end user, for high quality experiences that stay on-brand.

The **model layer** is your foundational model, including any fine-tuning performed by you or the model developer to improve performance and safety.



# Azure AI Content Safety

Keep your AI content safe with tailored protection, real-time threat detection, and customizable controls.



## Multimodal filtering

Scan text, images and multi-media to identify, block and monitor harmful content



## Customized systems

Create blocklists and custom categories to block entire topics, not just specific words



## Prompt shielding

Identify and mitigate prompts that could expose you to prompt injection attacks



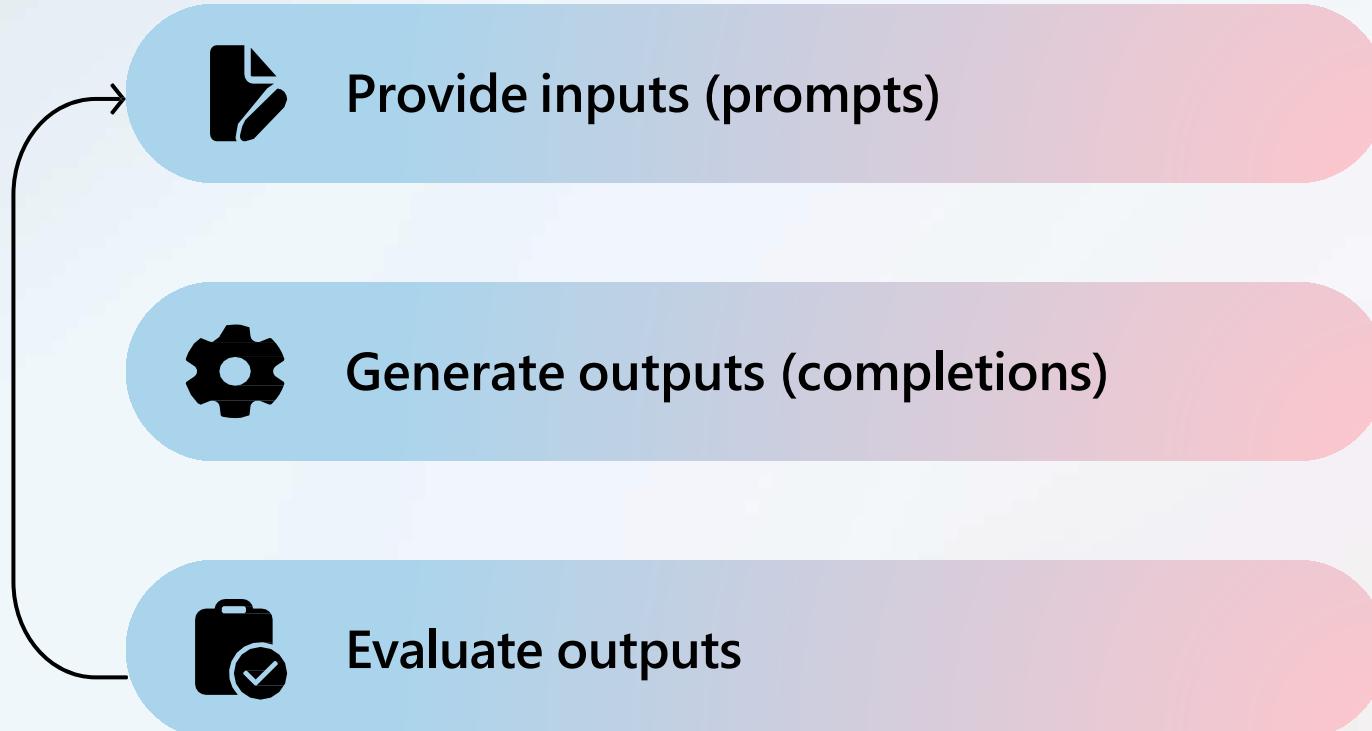
## Protected materials

Avoid outputting known or owned text content with protected materials detection



# Evaluations in a nutshell

**Learnings**



# Evaluations

## Quality

Groundedness

Coherence

Fluency

Relevance

Retrieval Score

Similarity

NLP Metrics (e.g., F1 Score)

## Risk & Safety

Jailbreak Defect

Hate and Unfairness

Sexual

Violence

Self-Harm

Protected Material

## Agents

Intent Resolution

Tool Call Accuracy

Task Adherence

Response Completeness

## Risk & Safety

Ungrounded Attributes

Code Vulnerability

# Evaluation best practices



Multiply your datasets and metrics



Evaluate systematically using GenAIOps



Review results as a team, with multiple points of view

# Getting started

# Accelerate AI deployment with Azure AI Foundry

[ai.azure.com](https://ai.azure.com)



Define scope & requirements



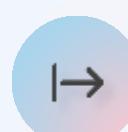
Prepare data



Build, train, and evaluate



Review and approve app for production



Deploy and monitor



Govern and monitor continuously

# Thank You!



Lino Tadros  
Co-founder & Chief Executive Officer

 lino@tahubu.com

 <https://www.linkedin.com/in/linotadros/>

