



# Microsoft Azure AI Services: Introduction

Martin Abrle, Partner Solution Architect

Juan Manuel Servera, Sr. Partner Solution Architect

Partner Organization, Microsoft Switzerland

# Agenda

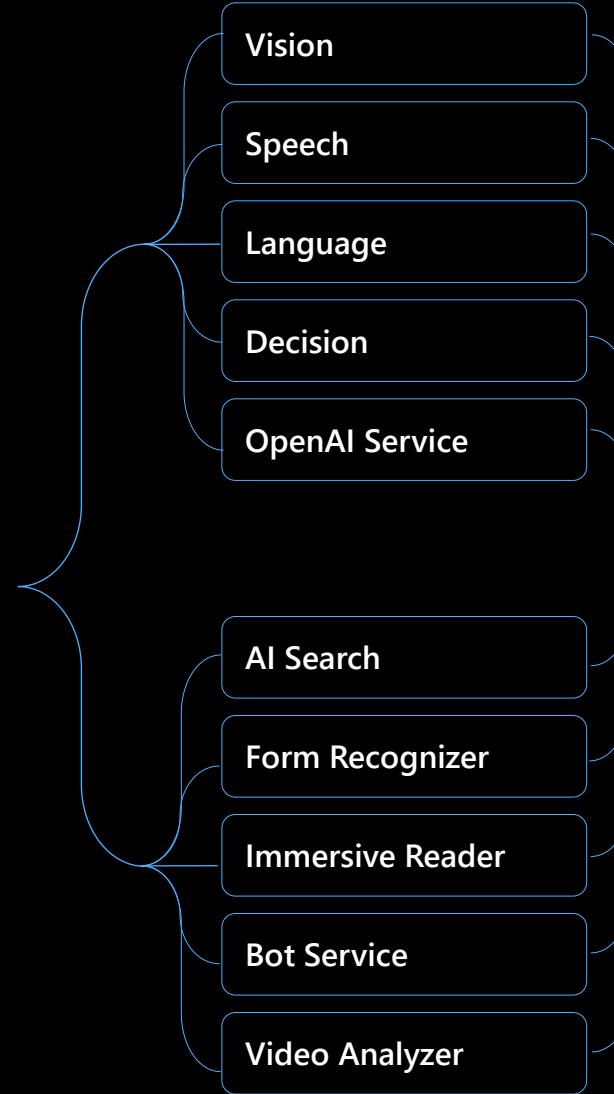
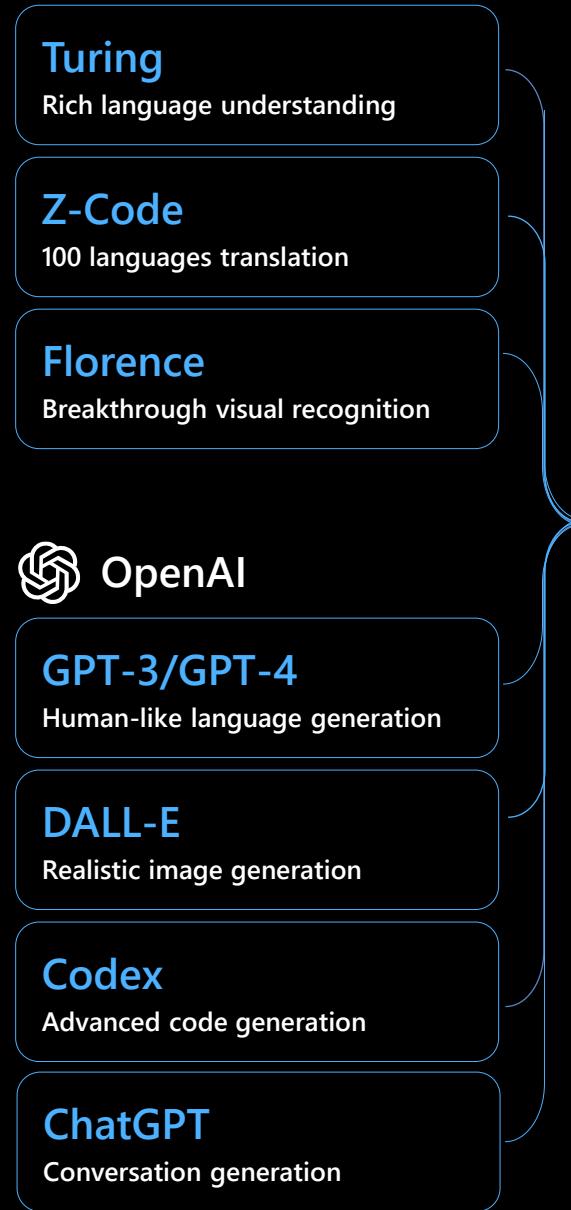
Cognitive services

Azure Open AI

AI Discovery Cards  
(afternoon)

# Azure AI Services

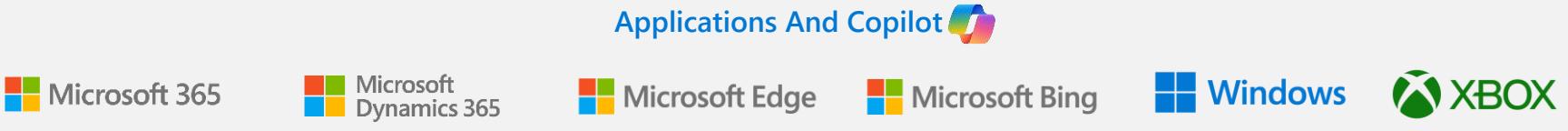
# Azure AI services



# Microsoft AI portfolio



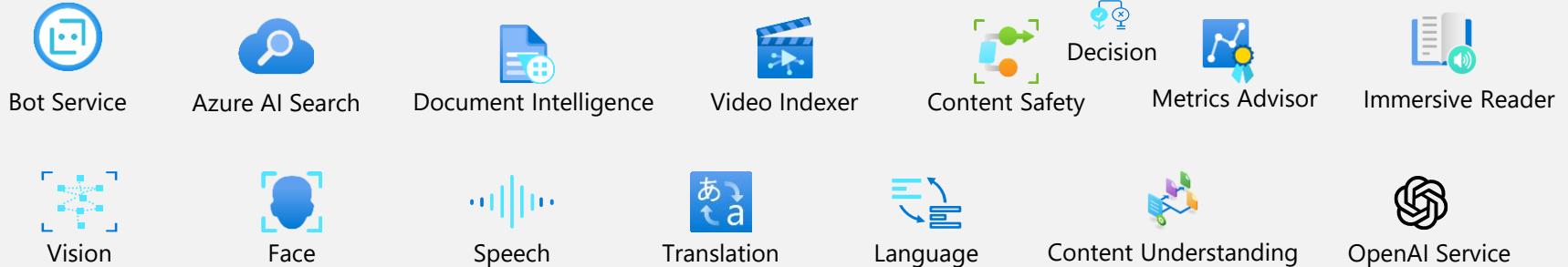
Business  
Users &  
Citizen  
Developers



Developers &  
Data Scientists

## Azure AI

### Azure AI services



### ML Platform



### Gen AI Platform



### Dev Platform



# Innovate responsibly

Microsoft improves facial recognition technology to perform well across all skin tones, genders

June 26, 2018 | John Roach

[in](#) [tw](#) [f](#) [re](#) [o](#)



## Responsible AI principles

Fairness



Reliability & Safety



Privacy & Security



Inclusiveness



Transparency



Accountability



Co-Designing Checklists to Understand Organizational Challenges and Opportunities around Fairness in AI

Michael Madan, Luke Stark, Jennifer Wotman Vaughan, Hanna Wallach  
CHI Conference on Human Factors in Computing Systems | March 2020  
Designated by ACM  
CHI 2020 Best Paper Award

[Download BibTeX](#)

[Download PDF](#)

**Groups**

FATE: Fairness, Accountability, Transparency, and Ethics in AI  
FATE | Montreal

Solving the challenge of securing AI and machine learning systems

Dec 6, 2019 | Valecia Maclin - General Manager, Engineering - Customer Security & Trust



Today, in collaboration with Harvard University's [Berkman Klein Center](#), we at Microsoft are publishing a series of materials we believe will contribute to solving a major challenge to securing artificial intelligence and machine learning systems. In short, there is no common terminology today to discuss security threats to these systems and methods to mitigate them, and we hope these new materials will provide baseline language that will enable the research community to better collaborate.

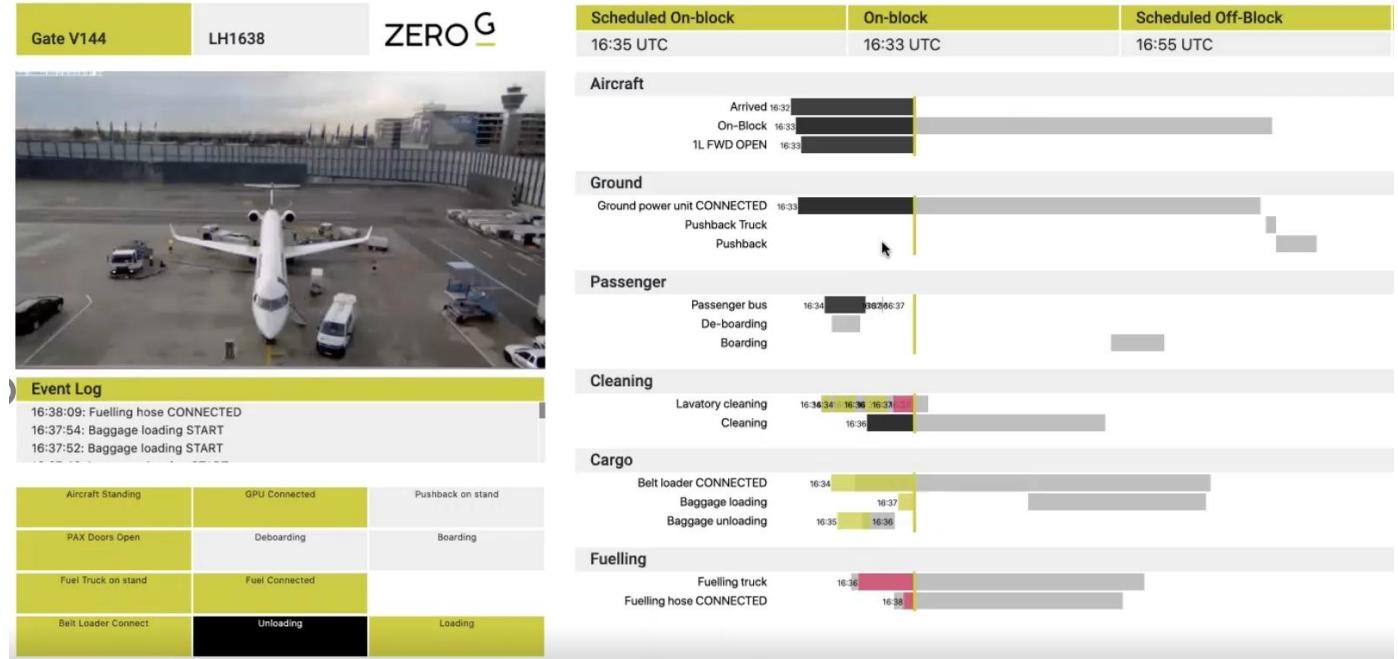


# Azure Video Analyzer

Gather and manage video insights

Tailor to your needs

Build in days, not months



## USE CASES



Process optimization



Workplace safety



Digital asset management

# Azure Metrics Advisor

Proactively monitor metrics and diagnose issues

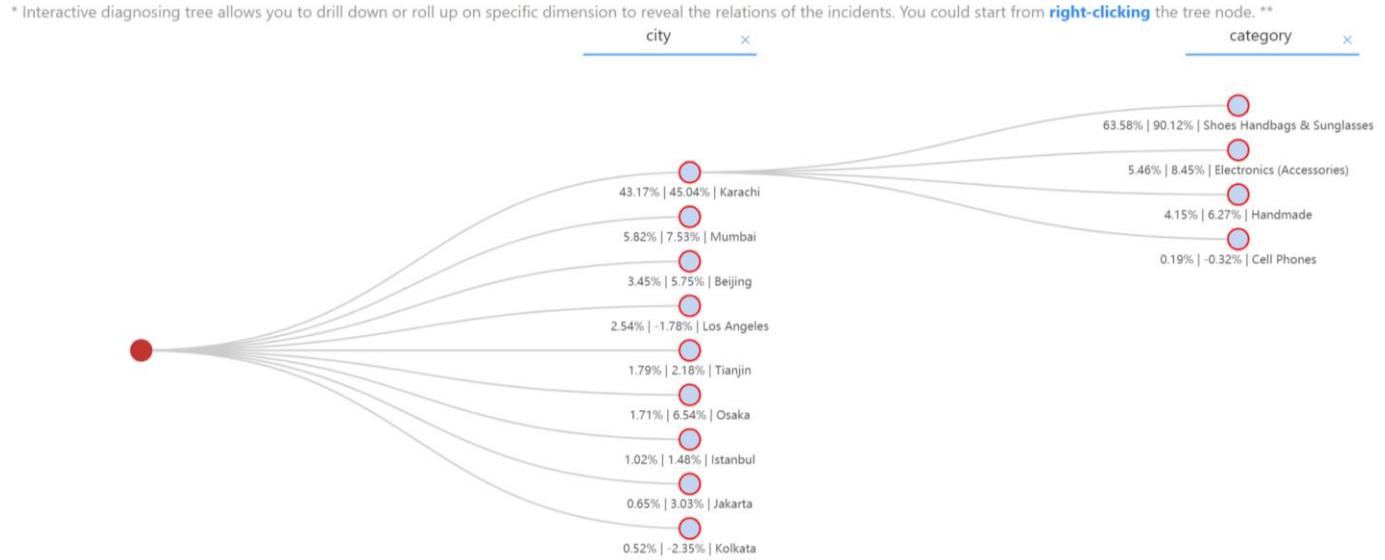
Detect anomalies for a variety of scenarios

Focus on fixing rather than monitoring

Accelerate time to value



Quick Diagnosing Tree   Interactive Tree   Metric Drilling   Similar Time-series Clustering   Compare Tools



## USE CASES



Business metrics monitoring

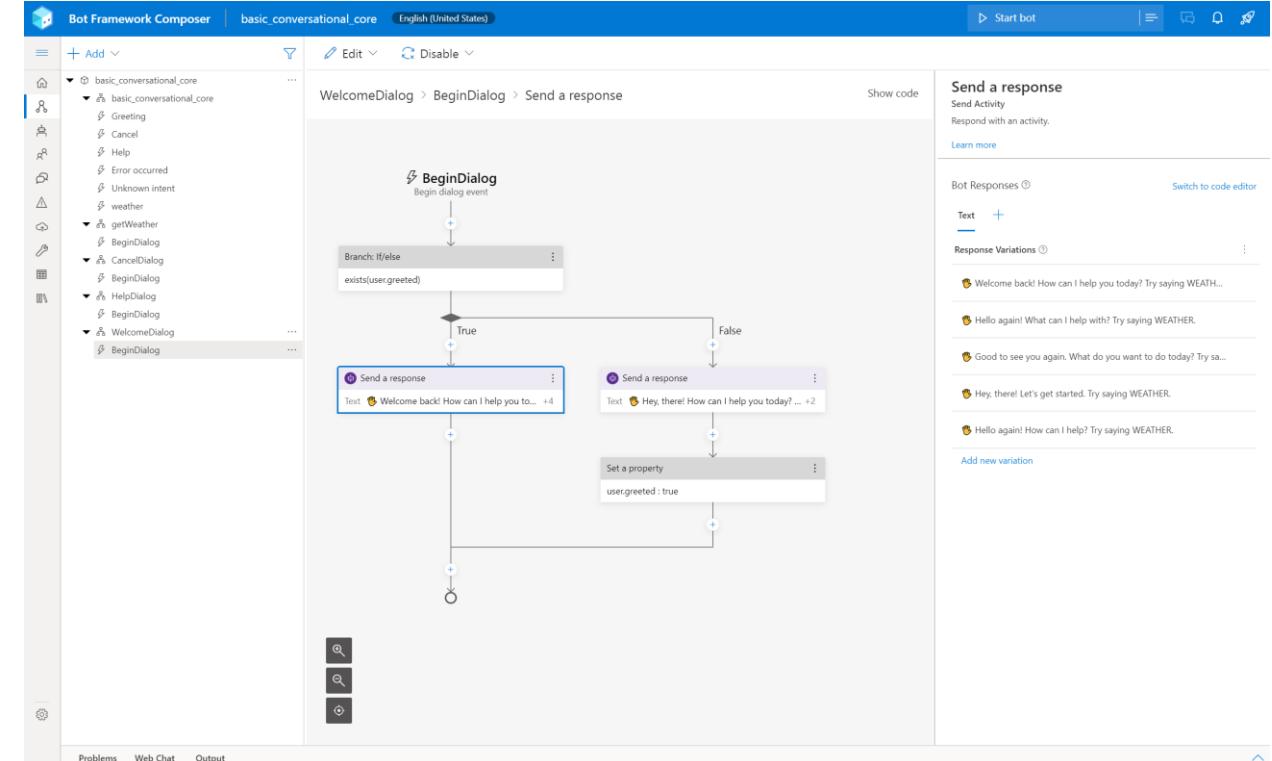
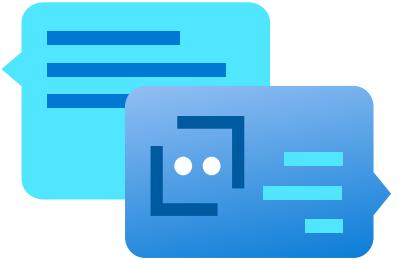


AIOps



Predictive maintenance

# Azure Bot Service



Accelerate bot building

Deploy to channels with minimal code changes

Conversational AI for healthcare

## USE CASES



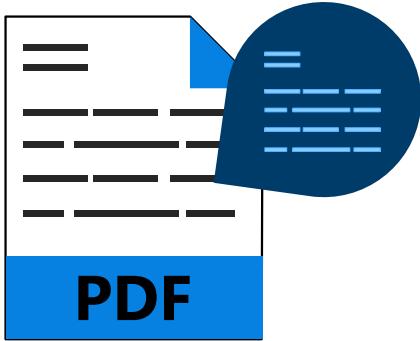
Customer service  
assistant



Smart  
personalized care



Enable voice control  
(voice assistant)



# Azure AI Document Intelligence

Extract text and structure with simple API

Customize extraction to your forms

Deploy anywhere, from cloud to edge

Export to Markdown

The screenshot shows the Azure AI Document Intelligence interface. On the left is a sidebar with various icons for file operations. The main area displays a PDF invoice from 'CONTOSO LTD.' to 'Microsoft Corp'. The invoice details include the date range from 07/02/2019 to 07/07/2019, service period from 07/02/2019 to 07/07/2019, and customer ID 'MS-12345'. The results panel on the right lists extracted entities with their confidence scores:

Entity	Confidence Score
BillingAddressRecipient	95.40%
Microsoft Finance	95.10%
CustomerAddress	95.10%
123 Other St, Redmond WA, 98052	95.40%
CustomerAddressRecipient	95.40%
Microsoft Corp	96.10%
CustomerId	96.10%
CID-12345	94.60%
CustomerName	94.60%
MICROSOFT CORPORATION	96.90%
DueDate	96.90%
text: 12/15/2019	valueDate: 2019-12-15
InvoiceDate	96.70%
text: 11/15/2019	valueDate: 2019-11-15
InvoicelId	97.00%
INV-100	96.70%
InvoiceTotal	96.70%
text: \$110.00	valueNumber: 110
Items	95.60%
Click to view analyzed table	
PreviousUnpaidBalance	95.60%
text: \$500.00	valueNumber: 500
PurchaseOrder	96.20%
PO-3333	
RemittanceAddress	94.80%
123 Remit St New York, NY, 10001	
RemittanceAddressRecipient	95.40%
Contoso Billing	
ServiceAddress	94.60%

## USE CASES



Claims management  
and automation



Document process  
automation



Government forms  
automation



# Azure AI Search

Start, maintain, and scale with minimal investment

Create searchable content using integrated AI

Search with user intent and not just keywords

Hybrid Search for RAG

The screenshot shows the Azure AI Search portal interface. On the left is a navigation sidebar with sections like Overview, Activity log, Access control (IAM), Tags, Diagnose and solve problems, Settings (Knowledge Center, Keys, Scale, Search traffic analytics, Identity, Networking, Properties, Locks), Monitoring (Alerts, Metrics, Diagnostic settings, Logs), Automation (Tasks (preview), Export template), Support + troubleshooting (Resource health, New support request), and Essentials. The main area displays resource details for a search service named 'oil-gas'. It shows the Resource group (change) as 'oil-gas', Status as 'Running', Location as 'West US 2', Subscription (change) as 'CSI POC Subscription\_LuisCa', Subscription ID as '7968f047-21ab-4b66-a5f9-1608f039b7a9', and Tags (change) as 'Click here to add tags'. The URL is https://oil-gas-s1.search.windows.net. It also shows Pricing tier as 'Standard', Replicas as '1 (No SLA)', Partitions as '1', and Search units as '1'. Below this, there's a section for semantic search capabilities with a 'Get Started' button and a link to 'View Cost' and 'JSON View'. The bottom right features a summary of AI search capabilities: 'Connect your data' (Import Data, Learn more), 'Use AI to extract and enrich' (Learn more), 'Explore your data' (Launch Explorer, Learn more), and 'Monitoring and administration' (Optimize Performance, Learn more).

## USE CASES



Content search



Product discovery optimization



Digital asset management



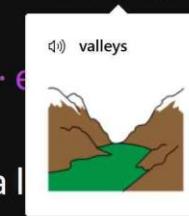
# Azure Immersive Reader

Empower users of all ages and abilities

Get started quickly



The study of Earth's landforms is called physical geography. Landforms can be mountains and valleys. They can also be glaciers or rivers. Landforms are sometimes called physical features. It is important for students to know about the physical geography of Earth. The seasons, the



## USE CASES



Enhance reading comprehension



Multilingual learning



Accessible learning

# Azure Immersive Reader



# State-of-the-art Computer Vision Features

Extract robust insights from images and video content across multiple industry domains

## Manage Digital Assets



Image  
Retrieval



Automatic,  
dense captions



Background  
Removal

## Enhance Safety & Security



Video  
Search



Video  
summarization

## Automate Retail Operations



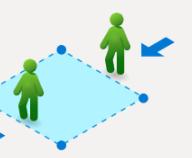
Model  
Customization



Product  
Recognition

Try it Out at  
**<https://aka.ms/VisionStudio>**

# Azure Vision AI

 <p><b>Extract common tags from images</b></p> <p>Use an AI model to automatically assign one or more labels to an image.</p> <p><a href="#">Try it out</a></p>	 <p><b>Detect common objects in images</b></p> <p>Recognize the location of objects of interest in an image and assign them a label.</p> <p><a href="#">Try it out</a></p>	 <p><b>Extract text from images</b></p> <p>Extract printed and handwritten style text from images and documents for supported languages.</p> <p><a href="#">Try it out</a></p>	 <p><b>Search photos with natural language</b></p> <p>Preview</p> <p>Retrieve specific moments within your photo album. For example, you can search for: a wedding you attended last summer, your pet, or your favorite city.</p> <p><a href="#">Try it out</a></p>	 <p><b>Add captions to images</b></p> <p>Preview</p> <p>Generate a human-readable sentence that describes the content of an image.</p> <p><a href="#">Try it out</a></p>	 <p><b>Dense captioning</b></p> <p>Preview</p> <p>Generate human-readable captions for all important objects detected in your image</p> <p><a href="#">Try it out</a></p>
 <p><b>Customize models</b></p> <p>Train custom image classification and object detection models using Vision Studio and Azure ML.</p> <p><a href="#">Start a project</a></p>	 <p><b>Video summary and frame locator</b></p> <p>Preview</p> <p>Generate a brief summary of the main points shown in video. Locate specific keywords and jump to the relevant section.</p> <p><a href="#">Try it out</a></p>	 <p><b>Scan cameras for defects</b></p> <p>Scan cameras for defects like blurry, under and over exposed cameras.</p> <p><a href="#">Learn more</a></p>	 <p><b>Count people in an area</b></p> <p>Analyze real-time video to count the number of people in a designated zone in a camera's field of view.</p> <p><a href="#">Learn more</a></p>	 <p><b>Detect when people cross a line</b></p> <p>Analyze real-time streaming video to detect when a person crosses a line in the camera's field of view.</p> <p><a href="#">Learn more</a></p>	 <p><b>Detect when people enter/exit a zone</b></p> <p>Analyze real-time streaming video to detect when a person enters or exits a zone in the camera's field of view.</p> <p><a href="#">Learn more</a></p>

Video Analysis

# Image Retrieval APIs

Vision Studio > Search photos with natural language

## Search photos with natural language PREVIEW

Retrieve specific moments within your photo album. For example, you can query: a wedding you attended last summer, your pet, your favorite city. Search for images based on the content of the image itself, rather than relying solely on manually assigned keywords or tags

Platforms

Cloud

[View documentation](#) [View SDK reference](#) [Use the REST API](#) [View samples on Github](#)

### Try it out

To try out this feature choose from a sample below. To try searching your own images, [sign in with Azure](#)

Sample image sets

Try with your own images



Nature

No.of photos: 260



Manufacturing

No.of photos: 245



Education

No.of photos: 264



Retail

No.of photos: 265

### Select a retrieval query or create your own

Enter a custom query

[Search](#) [Reset search](#)

### Query results

Query results vary from most relevant in the dataset to least relevant. Utilize the slider below to view more or less images based on their relevance to the retrieval query.

Most relevant



Least relevant



# Image Analysis APIs

## REST API

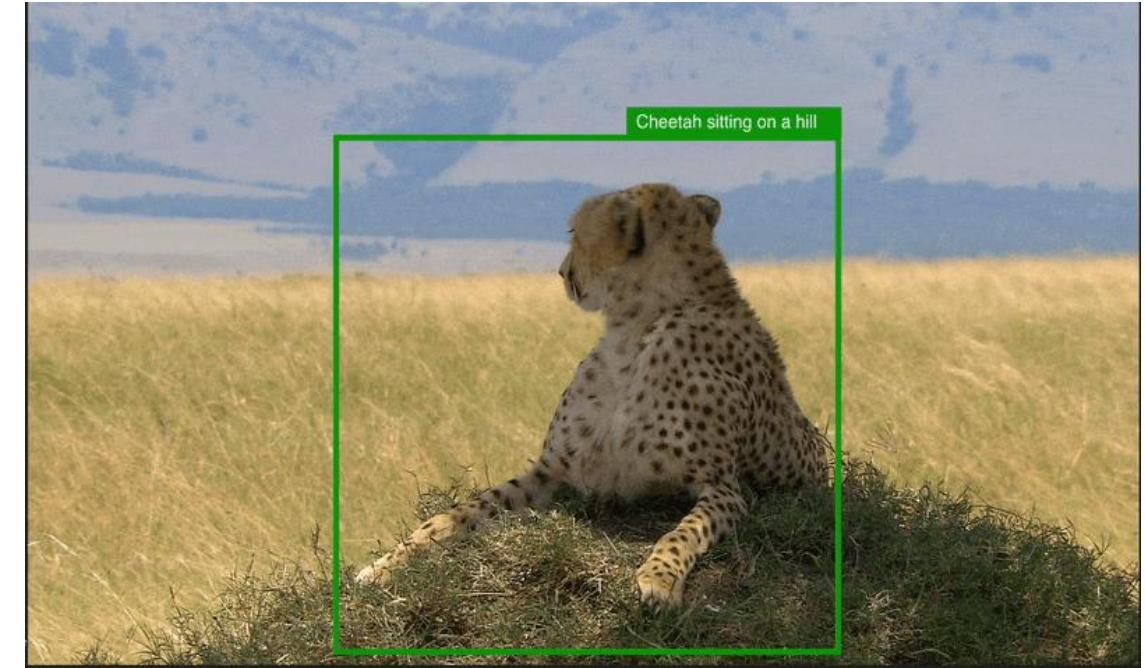
<https://{{endpoint}}/computervision/imageanalysis:analyze?features=caption>

<https://{{endpoint}}/computervision/imageanalysis:analyze?features=densecaptions>



a woman in a canoe touching water with a man in the front

Image captions



Dense captions

# Segment API: Background Removal

REST API

`https://{{endpoint}}/computervision/imageanalysis:segment?mode=backgroundremoval`



## Video summarization and frame locator PREVIEW



Video search and summarization uses a combination of natural language processing and computer vision techniques to analyze the content of a video. It can quickly and concisely summarize the main points of a video, and it also allows you to search for specific moments within the video, making it easy to find relevant content.

Platforms

Cloud

Interested in trying out the APIs? Apply [here](#) for access.

### Try it out

Choose a video clip to see the summarization and frame locator capabilities.

Note: videos that have been uploaded to Vision Studio will be stored in your account for 48 hours for this try out experience, after which they will be deleted automatically.

Drag and drop a file here  
or  
Browse for a file  
or  
Browse container for your video



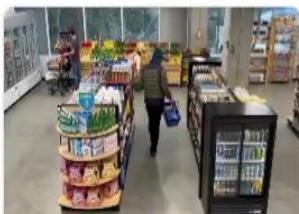
Datacenter

A video that showcases a data center, depicting the hardware that powers it and the people that maintain it.



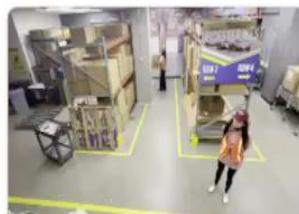
Factory

A video that depicts unsafe actions in factory, such as running, carrying boxes with hands, taking off helmets, etc.



Retail

A grocery store where people go to buy food and day-to-day items.



Warehouse

A warehouse where the safety of the workers is a top priority.

### Run a test

Choose a test

Run

# Video Retrieval

# Model customization

**Fast path for customization for specific use cases**

Few Shot learning and  
model bootstrapping with  
single image per label

Active Learning Loop with  
production data

Continuous learning for  
new use cases



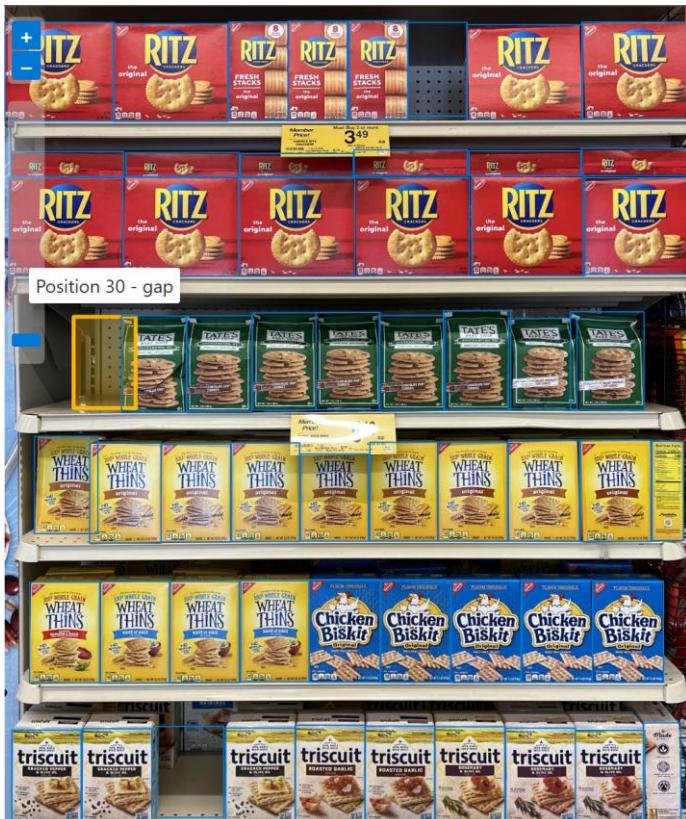
Custom Model

# Product Recognition APIs

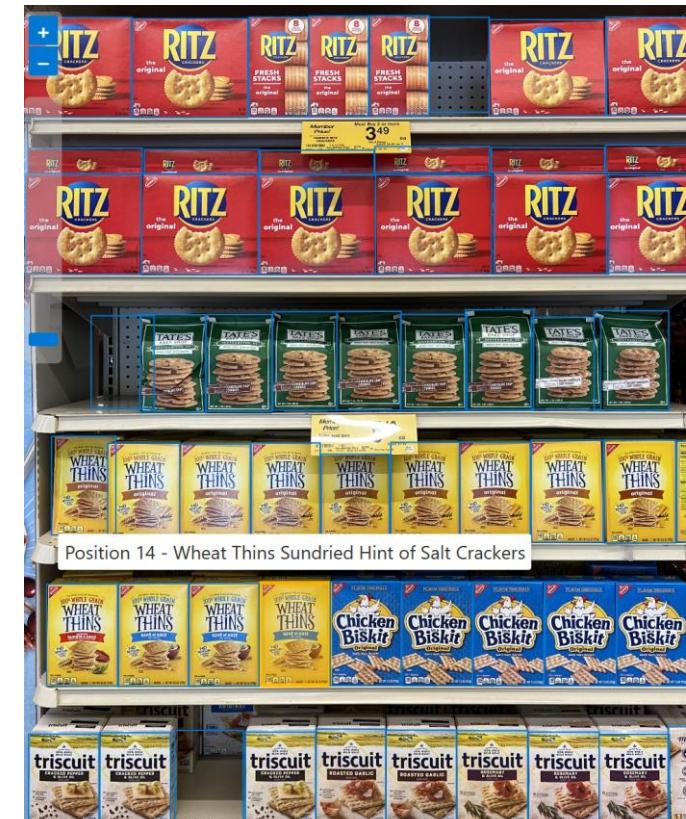
## Product Recognition

<https://{{endpoint}}/computervision/operations/shelfanalysis-productunderstanding:analyze>

<https://{{endpoint}}/computervision/operations/shelfanalysis-productunderstanding:analyze?product-classifier-model-name=custommodel>



\*



\*

**Demos:  
No code, Try out experience available**

Vision Studio – <https://aka.ms/visionstudio>

# Azure Cognitive Services

## Speech Service



### Speech to Text

Language Comprehension

Speech Transcription

Versatile Speech Models



### Text to Speech

Neural Voices

Custom Neural Voices

Cross lingual  
Adaptation



### Speech Translation

Language Detection

Text-to-Text Translation

Speech-to-Text  
Translation



### Speaker ID

Speaker Identification

Speaker Verification

Diarizations

SPEECH SDK

DEVICES SDK

AI CONTAINERS

LANGUAGES

CUSTOM MODELS

# Batch versus Real Time Playbook

## How does one decide on the operation mode

Understand the latency requirements of your solution

Understand the cost delta (real time has higher set costs)

Understand development support requirements (You are likely to require development resources)

### Batch

Faster to set up through V3 REST API

More features (diarization)

Lower response times (a 10min file is returned within 5mins + queueing time)

### Real Time

Requires ingress management (SDK instances on VMs at a size relative to the audio volume)

Faster than real time for file streaming (a 10 min file is transcribed in 5mins)

Higher maintenance costs as well as up front development costs

# Custom neural voices

## Brand identity

Design & implement a voice model that complements or augments your brand strategy

## Custom persona

Extract value from your analytics to work by accommodating customer sentiment with custom voice characteristics

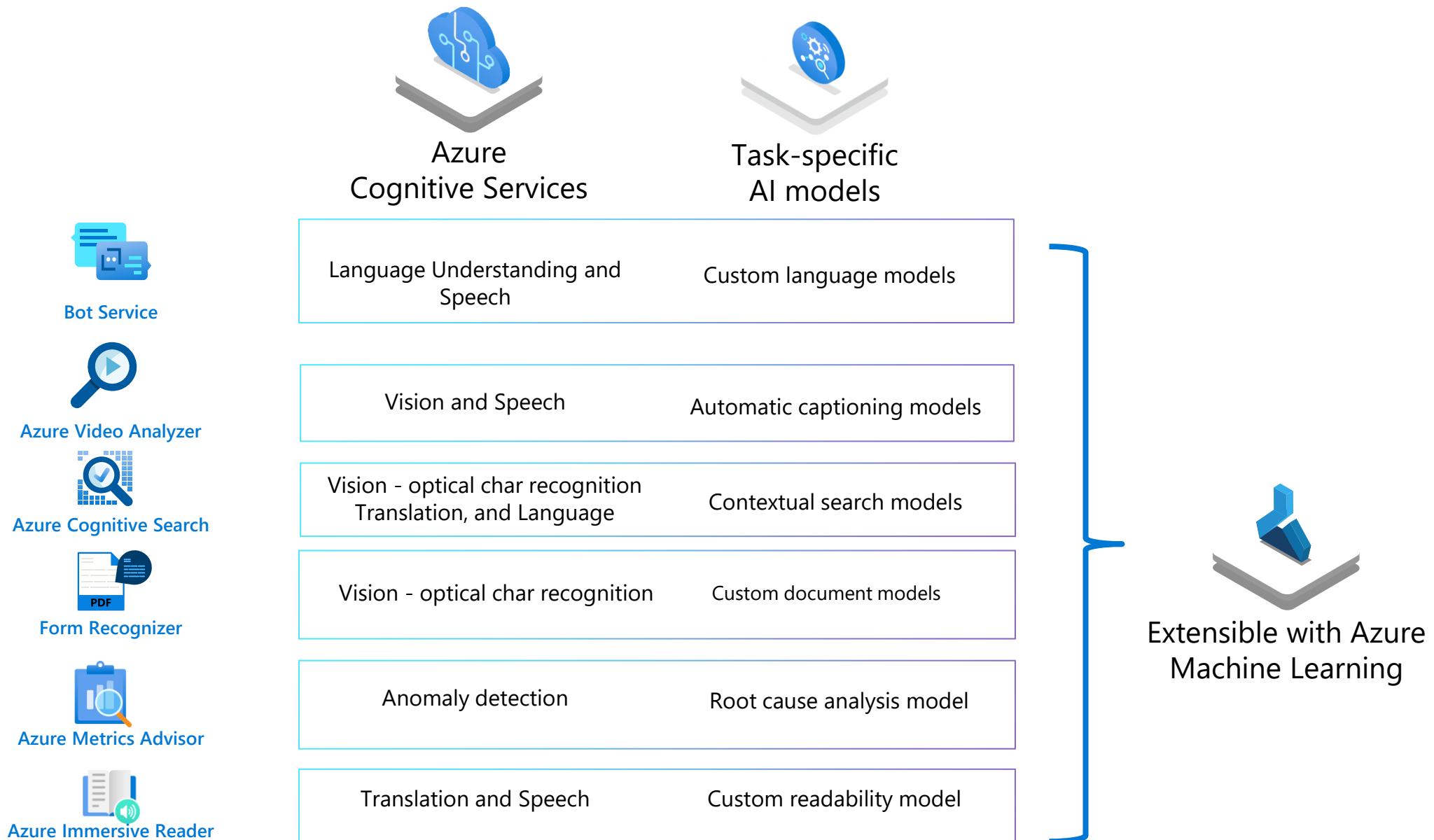
## Natural interaction

Increase your customer's emotional connection and interaction with your applications

The screenshot shows the Microsoft Cognitive Services | Speech Studio portal at <https://speech.microsoft.com/portal>. The page features a survey link, a text-to-speech overview, and three main sections: 'Voice Gallery', 'Custom Voice', and 'Audio Content'.

- Voice Gallery:** Quickly pick up a voice for your text-to-speech apps. Includes a 'Try it out' button.
- Custom Voice:** Build a recognizable one-of-a-kind voice for your text-to-speech apps with your available speaking data. Includes a 'Start a project' button.
- Audio Content:** Visually control A, such as voice style and breaks. Quick and customized a. Includes a 'Start a project' button.

# Modernize business processes



# Azure OpenAI

# Azure OpenAI

*Azure OpenAI Service is a REST API service, which provides access to OpenAI's powerful language models including o1, o1-mini, GPT-4o, GPT-4o mini, GPT-4 Turbo with Vision, GPT-4, GPT-3.5-Turbo, and Embeddings model series. These models can be easily adapted to your specific task including but not limited to content generation, summarization, image understanding, semantic search, and natural language to code translation. Users can access the service through REST APIs, Python SDK, or in the [Azure AI Foundry](#). ([link](#))*

# 60,000+ organizations using Azure AI today





# Azure OpenAI Service

99.9%

28

99%

EU & US

Reliability  
service wide

Regions  
for data residency

Latency SLA  
for Provisioned

Data Zones  
for regional control

# Using Azure OpenAI Service

**When you use Azure OpenAI Service, your prompts (inputs) and completions (outputs), your embeddings, and your training data**

Are **NOT** available to other customers.

ARE **NOT** available to OpenAI.

Are **NOT** used to improve OpenAI models.

Are **NOT** used to improve any Microsoft or 3<sup>rd</sup> party products or services.

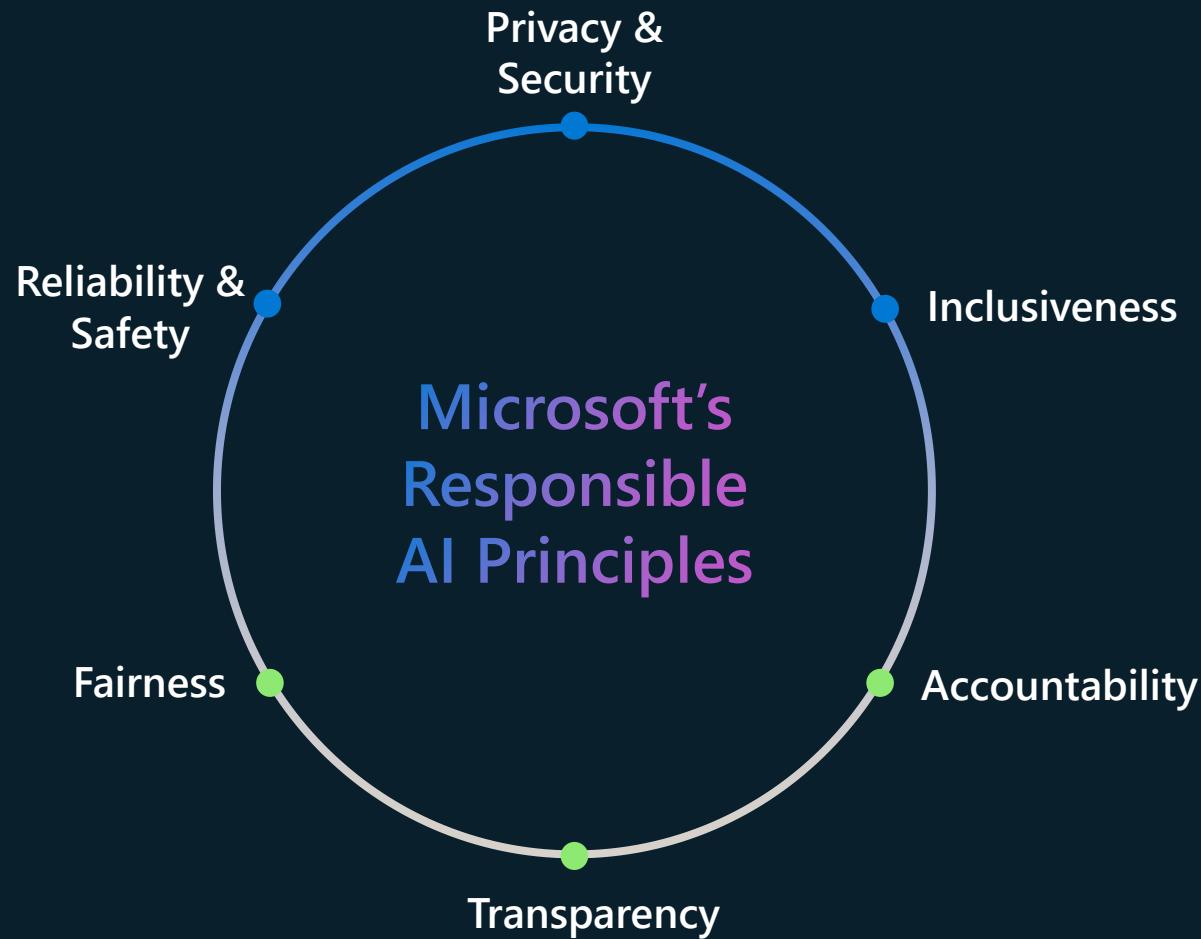
Are **NOT** used for automatically improving Azure OpenAI models for your use in your resource (The models are stateless, unless you explicitly fine-tune models with your training data).

Your fine-tuned Azure OpenAI models are available **exclusively** for your use.

The Azure OpenAI Service is fully controlled by Microsoft.

Microsoft hosts the OpenAI models in Microsoft's Azure environment and the Service does **NOT** interact with any services operated by OpenAI (e.g., ChatGPT, or the OpenAI API).

# Responsible AI protection turned on by default



# Azure AI Content Safety

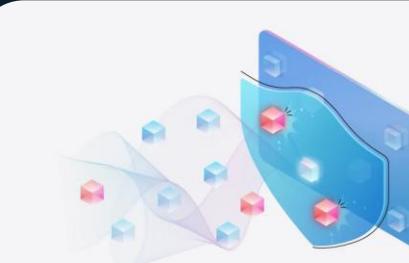
Safety is built-in with content filters for Azure OpenAI Service



## Harms, Fairness & Bias

Content filters for inputs + outputs:

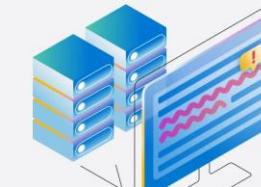
- Violence
- Hate and Fairness
- Sexual
- Self-Harm
- Custom categories PREVIEW
- Custom blocklists



## Security

Prompt shields for inputs:

- Direct prompt injection attacks
- Indirect prompt injection attacks

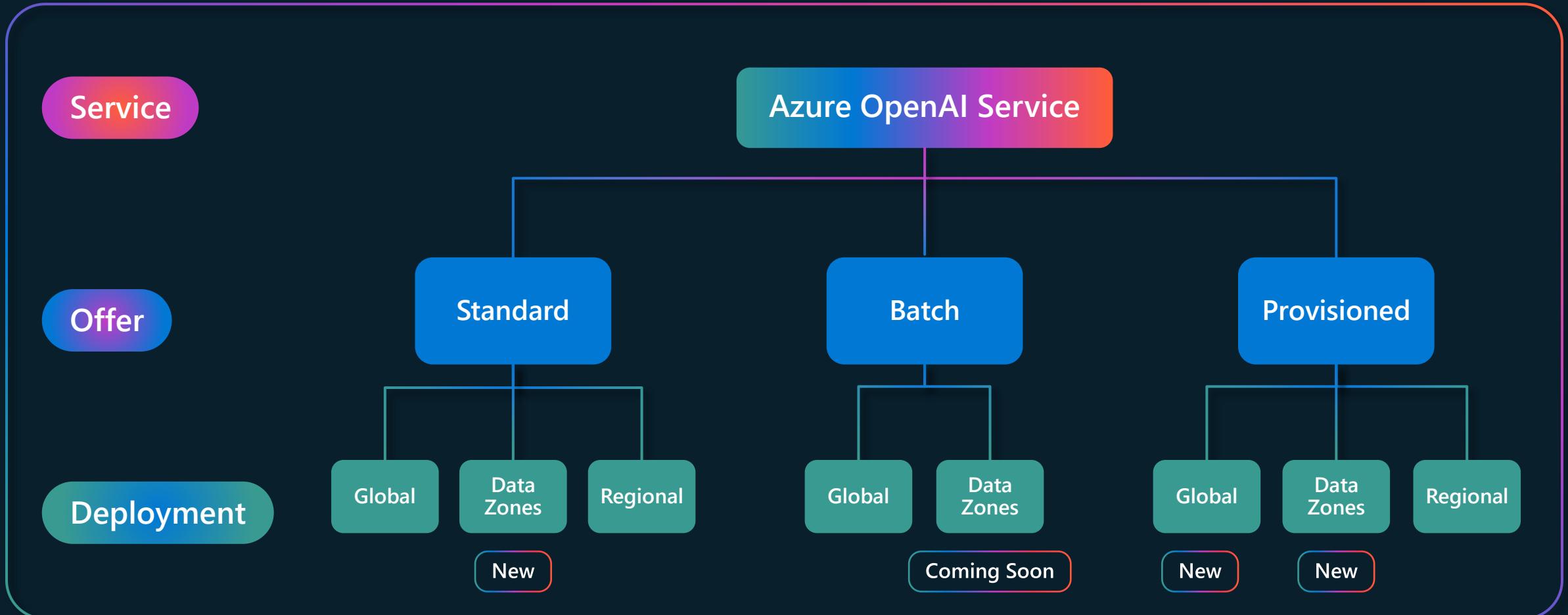


## Quality & Transparency

Content filters for inputs + outputs:

- Protected text
- Protected code PREVIEW

# Azure OpenAI Service Offerings



# Choosing the best offer

## Standard

### Benefits

- Flexible to scale up and down upon demand
- Optimized for low to medium volume workloads
- Ease of management

### RECOMMENDED FOR

- ✓ Production workloads
- ✓ Development and testing
- ✓ Prototyping and proof of concept

## Batch

### Benefits

- Efficient Processing of large volumes of data
- Scalable to handle varying workloads
- Cost-effective for large-scale deployments

### RECOMMENDED FOR

- ✓ Large-scale data processing
- ✓ Generate large volumes of content Transforming data at scale
- ✓ Evaluate LLM models and assess

## Provisioned

### Benefits

- Easy access with high and predictable throughput
- Real-time scoring for large consistent volume
- Cost effective for large-scale deployments

### RECOMMENDED FOR

- ✓ Production workloads
- ✓ High volume data processing
- ✓ Throughput heavy workloads
- ✓ Real-time

# Choosing the best deployment

## Global

### Overview

- Lowest price with highest throughput
- Broadest model availability
- Broadest capacity availability

### BEST FOR

- ✓ Applications requiring a consistent experience across multiple regions
- ✓ Services needing to be available globally with low latency
- ✓ Cost savings is a priority
- ✓ Data residency independent of deployment

## Data zone

### Overview

- Cross region load balancing within a geographic boundary (US or EU)
- Broader model availability
- Broader capacity availability

### BEST FOR

- ✓ Applications needing more processing capacity with data residency required
- ✓ Cost saving with meeting compliance requirements
- ✓ Optimal access to latest AI models and innovations

## Regional

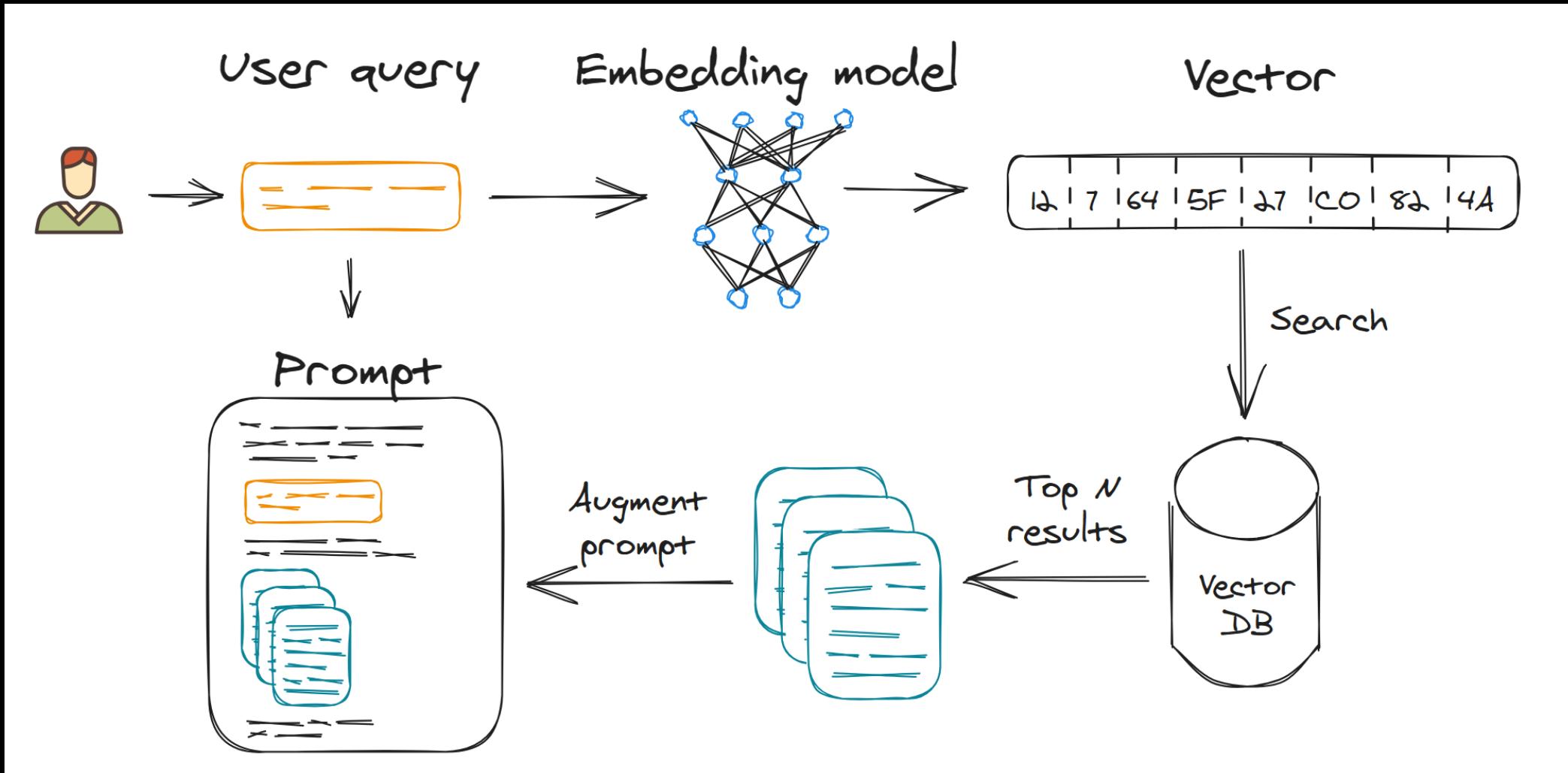
### Overview

- Strictest levels of data control
- Limited model availability
- Limited capacity availability

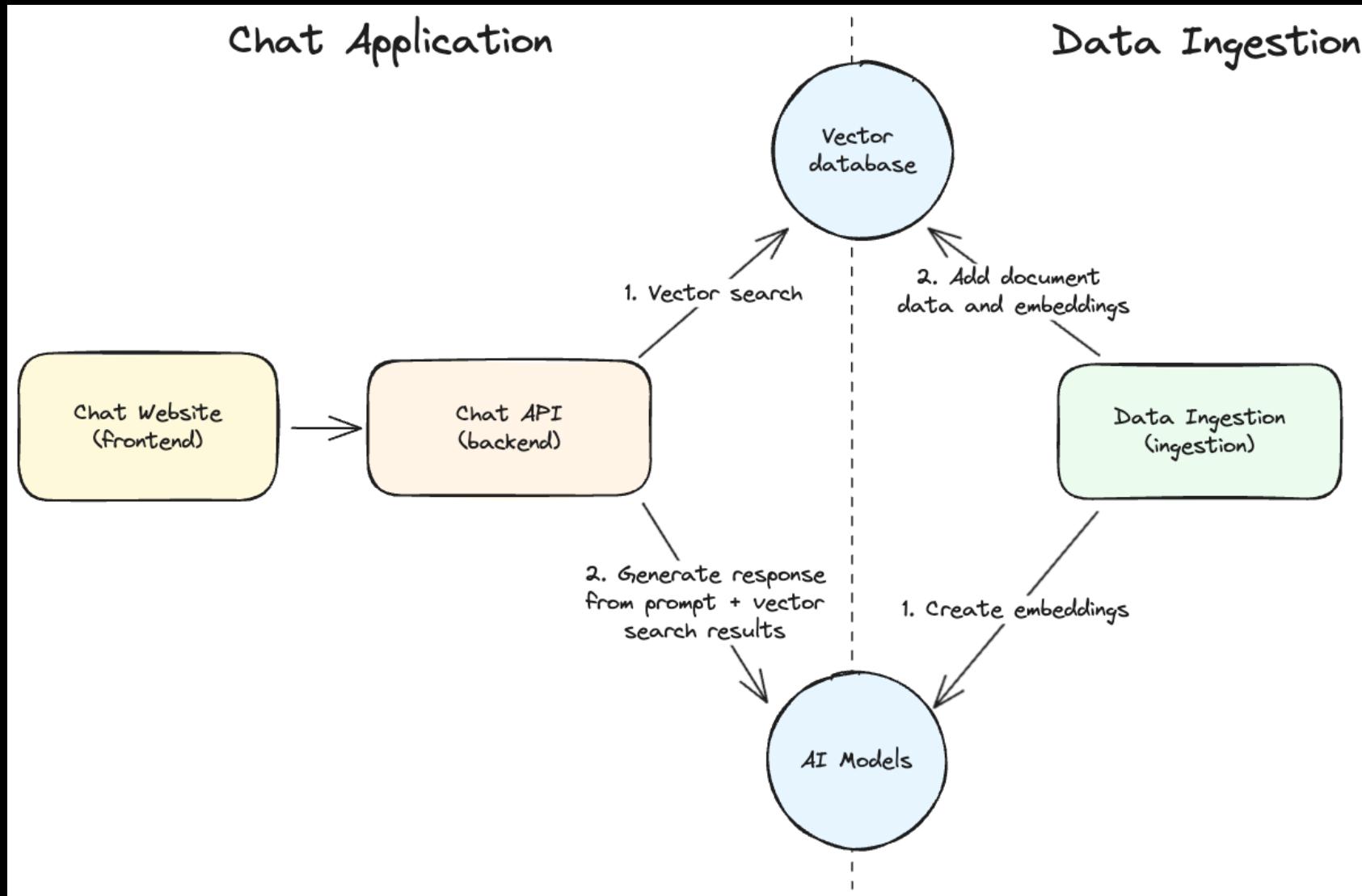
### BEST FOR

- ✓ Applications required to meet data residency compliance globally with low latency
- ✓ Services needing to be closer to the end-users to reduce latency
- ✓ Applications requiring localized data processing and storage

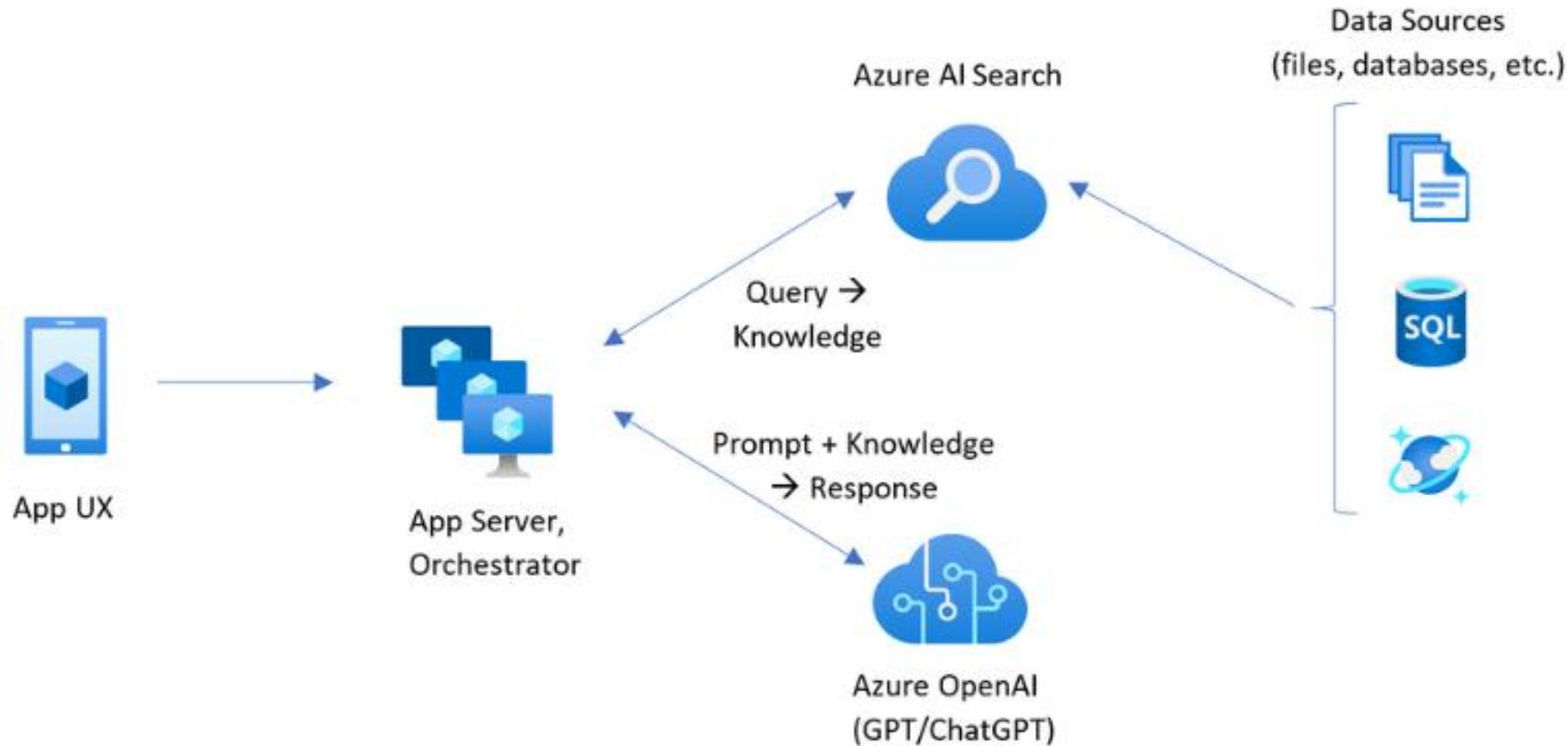
# RAG Application: user query processing



# RAG Application: architecture



# RAG Application on Azure: architecture ([link](#))



Microsoft runs on **Azure AI Foundry**

Microsoft  
Copilot

Microsoft 365  
Copilot

Microsoft Security  
Copilot

GitHub  
Copilot

DAX  
Copilot

# Demo

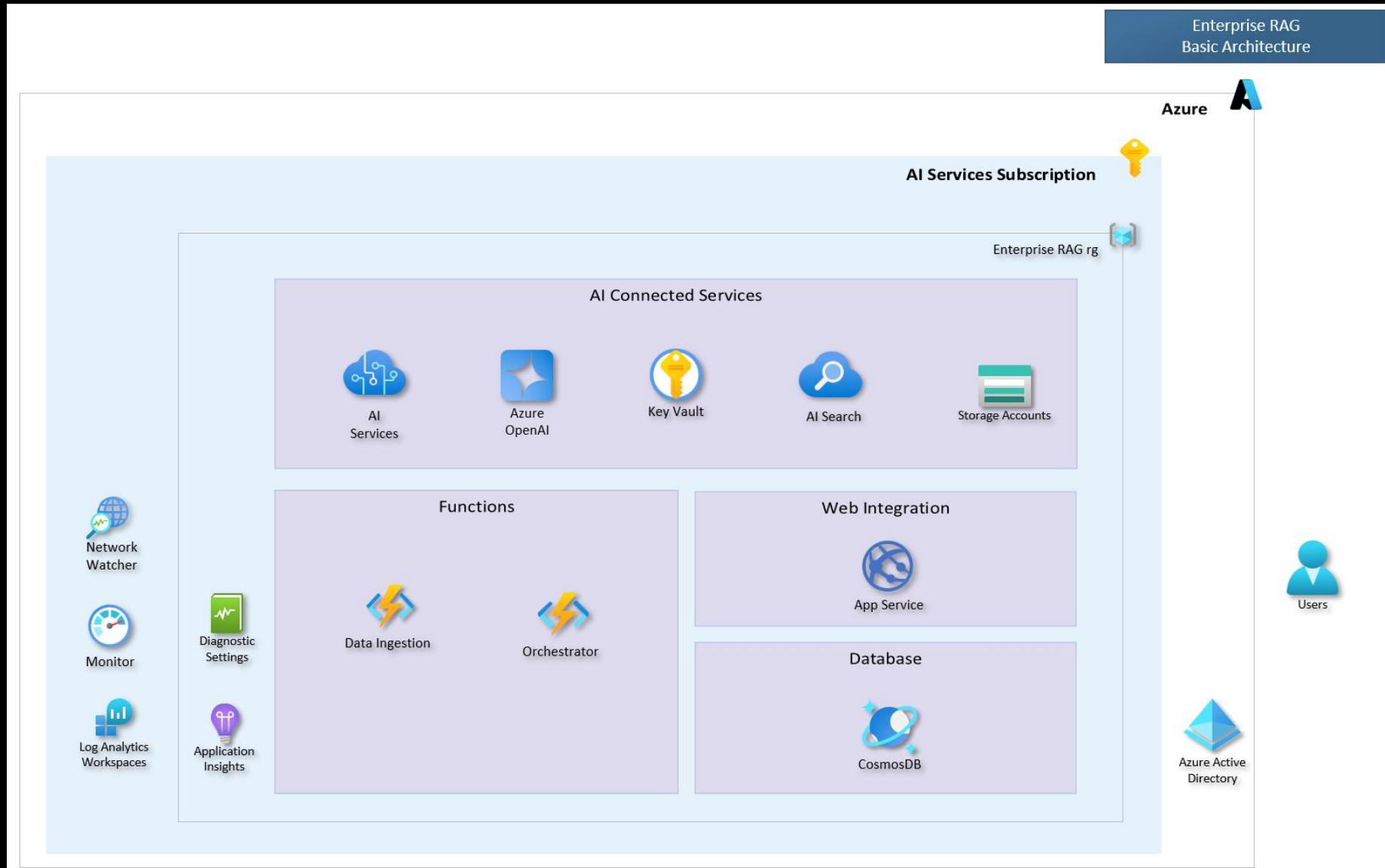
Azure Open  
AI

AI Foundry

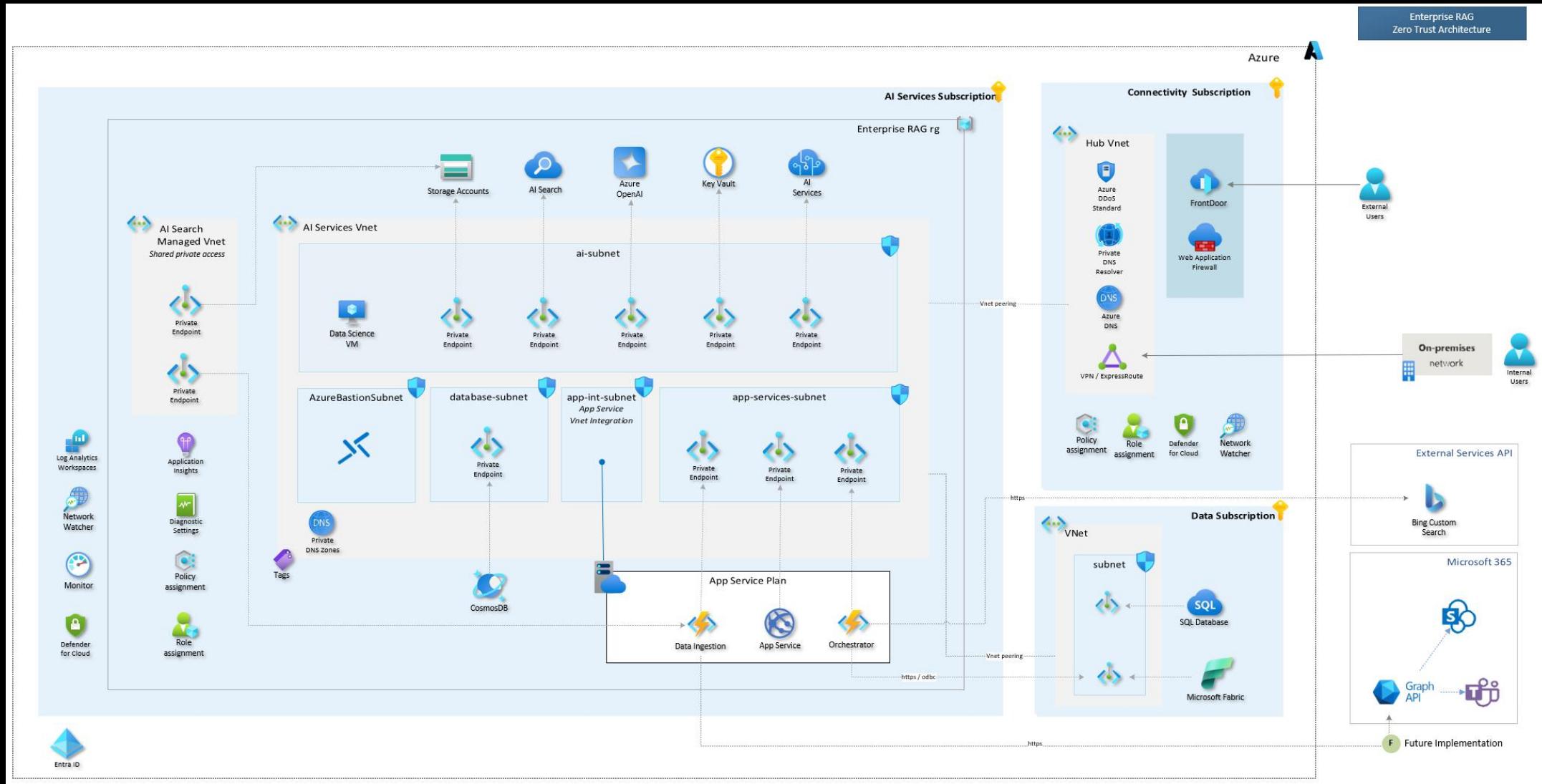
AI Search

RAG

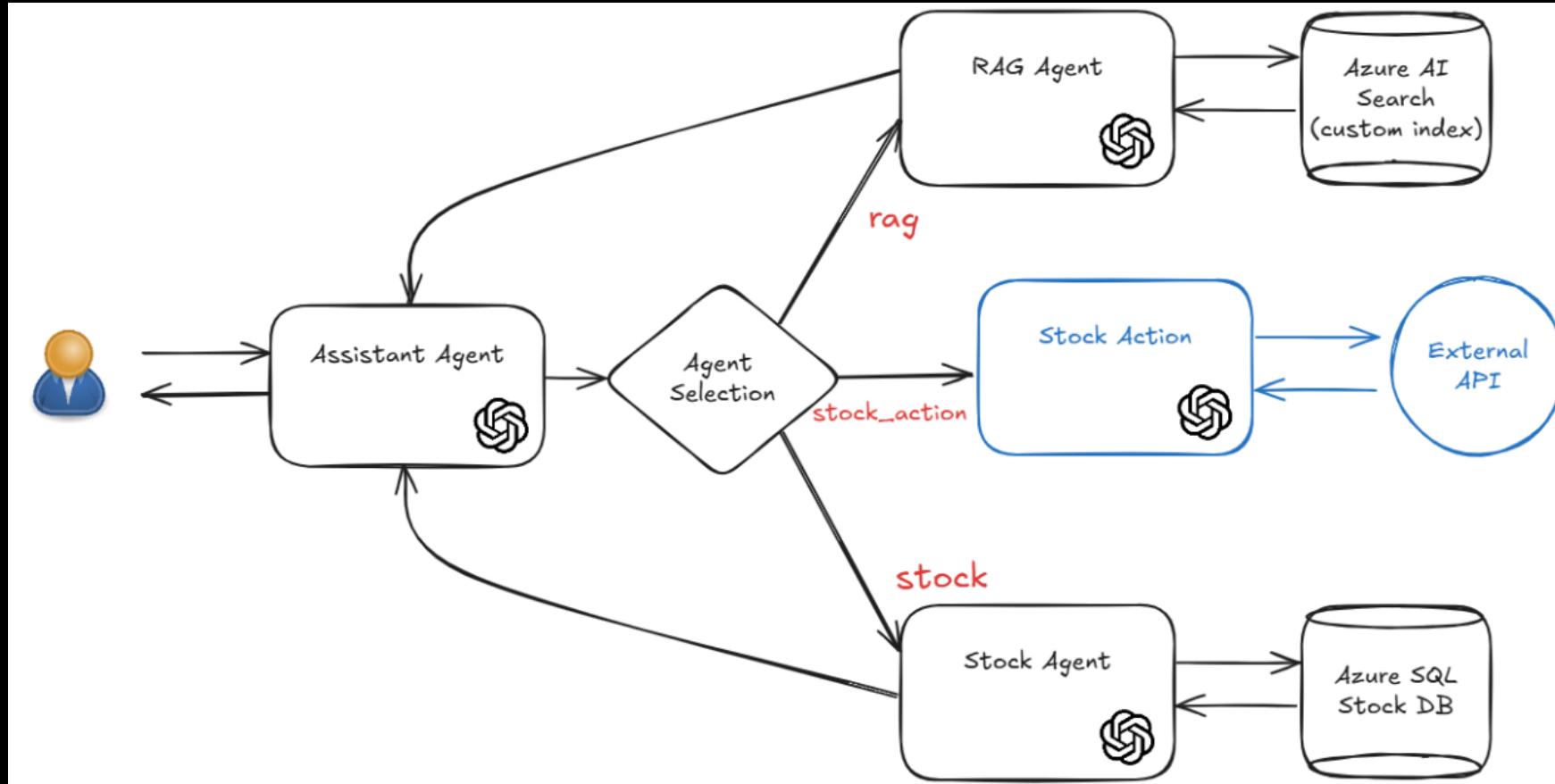
# RAG Application on Azure: more realistic architecture



# RAG Application on Azure: enterprise architecture



# Agentic workloads



Framework examples:

- Microsoft Autogen
- Phidata
- OpenAI Swarm
- CrewAI
- LangGraph
- ....

# phidata example

- [GH Link](#)



```
from phi.agent import Agent
from phi.model.openai import OpenAIChat
from phi.tools.duckduckgo import DuckDuckGo
from phi.tools.yfinance import YFinanceTools

web_agent = Agent(
    name="Web Agent",
    role="Search the web for information",
    model=OpenAIChat(id="gpt-4o"),
    tools=[DuckDuckGo()],
    instructions=["Always include sources"],
    show_tool_calls=True,
    markdown=True,
)

finance_agent = Agent(
    name="Finance Agent",
    role="Get financial data",
    model=OpenAIChat(id="gpt-4o"),
    tools=[YFinanceTools(stock_price=True, analyst_recommendations=True, company_info=True)],
    instructions=["Use tables to display data"],
    show_tool_calls=True,
    markdown=True,
)

agent_team = Agent(
    team=[web_agent, finance_agent],
    model=OpenAIChat(id="gpt-4o"),
    instructions=["Always include sources", "Use tables to display data"],
    show_tool_calls=True,
    markdown=True,
)

agent_team.print_response("Summarize analyst recommendations and share the latest news for NVDA", strea
```

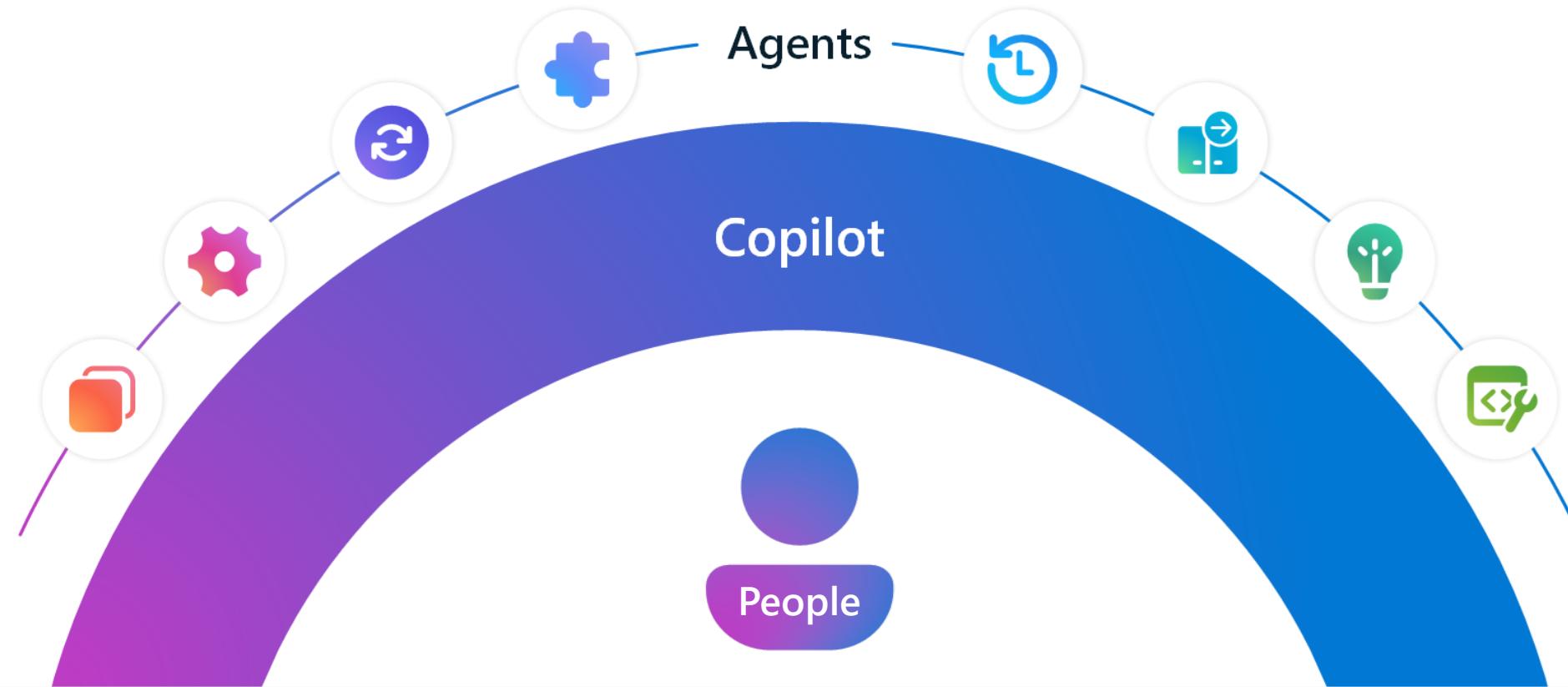


**Copilot**



**Business  
Process**

# Transform the way you work with Copilot + Agents



Copilot makes agents more personal to you

Agents make Copilot more capable

Together, they benefit from consistent security, governance & management

# Explore a continuum of solutions

## IT Helpdesk agent

How do I connect to the corporate network?



Simple



## Project Tracker agent

What is the status of phase 2 for project X and the remaining budget?



## Device Refresh agent

Request a new laptop and send approvals via IT Service tool.



## Budget Management agent

Review outstanding open PO's and begin financial planning.



## Lead Gen agent

The agent has identified and researched 15 new leads for you to review.



Advanced



## Customer Support agent

The agent has identified new support issues and triaged to other agents.



# Links

- <https://learn.microsoft.com/en-us/azure/ai-services/openai/how-to/provisioned-throughput-onboarding>
- <https://learn.microsoft.com/en-us/legal/cognitive-services/content-safety/data-privacy>
- <https://learn.microsoft.com/en-us/azure/ai-studio/how-to/concept-data-privacy>
- <https://learn.microsoft.com/en-us/azure/ai-services/openai/concepts/abuse-monitoring>



**Thank you**