# Agenda

# It's a crazy year for GenAI

NashTech has researched on popular generative AI models on LLMs, images and audio

**Azure OpenAI**
Jan 2023

**Copilot Studio**
**GPT-4**
Mar 2023

**Sematic Kernel**
**Promptflow**
May 2023

**Bing chat**
**Enterprise**
Jul 2023

**GPT-4 Vision**
Dec 2023

Mar 2023

**GPT-3.5**
**Copilot**
**GitHub Copilot X**

Apr 2023

**DALL.E 2**

Sep 2023

**Whisper**

Nov 2023

**Azure AI Studio**
**DALL.E 3**

Mar 2024

**Mistral**
**Llamma**

# Microsoft and OpenAI partnership

## OpenAI

Ensure that artificial general intelligence (AGI) benefits humanity

## Microsoft

Empower every person and organization on the planet to achieve more

## Azure OpenAI Service

| GPT-4, GPT-4 Turbo, GPT-3.5 Turbo | GPT-4 with Vision | Babbage & Davinci | DALL.E 3 | Whisper |
|---|---|---|---|---|
| Language | Multi-modal | Fine-tuning | Images | Transcription & Translation |

# DALL.E 3

**Azure OpenAI service**

DALL-E 3 is an image generation model that allows you to generate images from text prompts

*This image is generated by DALL-E 3*

# What can DALL.E 3 do?



**LOGO & BRANDING:** QUICK CONCEPT GENERATION.

**CREATIVE INSPIRATION:** OVERCOME DESIGN BLOCKS.

**CONTENT ILLUSTRATIONS:** UNIQUE IMAGES FOR BLOGS/ARTICLES.
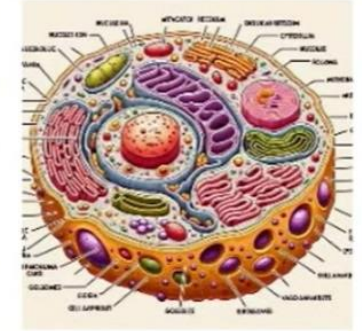
**AD CAMPAIGNS:** VISUALIZE MARKETING CONCEPTS.

**PRODUCT VISUALIZATION:** GAUGE INTEREST & FEEDBACK.

**EDUCATION:** CUSTOM IMAGERY FOR COURSES.

**FASHION DESIGN:** VISUALIZE CLOTHING PATTERNS.

**GAMING:** CHARACTER & ENVIRONMENT CONCEPTS.

# Announcing GPT-4V with Azure AI Vision

Unlock new scenarios with GPT-4V, Azure Open AI Service and Azure AI Vision integration

Add images to retrieval augment generation (RAG) patterns

Prompt with video, images, and text

# What GPT-4V offers?

**GPT-4 with Vision (GPT-4V)** is a multimodal model developed by OpenAI that accepts both image and text inputs and generates text outputs.

| Text prompt | + | Input image | = | Desired output text |

*Note: GPT-4V doesn't generate image outputs*

# Model catalog (in Azure AI & ML Studio)

**Catalog featuring the best as foundation model collections**

- Popular OSS models handpicked and optimized by AzureML

- Partnering with HuggingFace to offer thousands of OSS models for inference

- Azure OpenAI models

# MaaP and MaaS



| MaaP | | MaaS | |
|---|---|---|---|
| **User subscription** | | **User subscription** | |
| Inference endpoint | Fine-tuning and Evaluation jobs | Inference endpoint | Fine-tuning and Evaluation jobs |
| GPU | GPU | | |
| **Prod subscription** | | **Prod subscription** | |
| | | Model pool (compute) | |
| | | GPU | GPU |
| Model catalog (hosting) | | Model catalog (hosting) | |

| Total | $5 per hour | Consumption based offer | Provisioned offer |
|---|---|---|---|
| Software price (accrues to partner) | $2 per hour | **$0.01 / 1000 tokens** | **$2 per hour per scale unit** |
| Infrastructure price (accrues to Microsoft) | $3 per hour (with suggested VM instance) | Best effort latency and throughput | Guaranteed throughput: 1 scale unit supports 1000 requests per minute |

# Models as a Service

- Launched with Llama2 and Mistral

- Ready to use APIs with pay-as-you-go billing based on tokens for LLMs

- Integrate with your favorite LLM tools like PromptFlow, Semantic Kernel or LangChain

- Hosted finetuning without provisioning GPUs

# Provisioned Throughput
## Model processing capacity for your high-volume production workload

**Predictable performance**

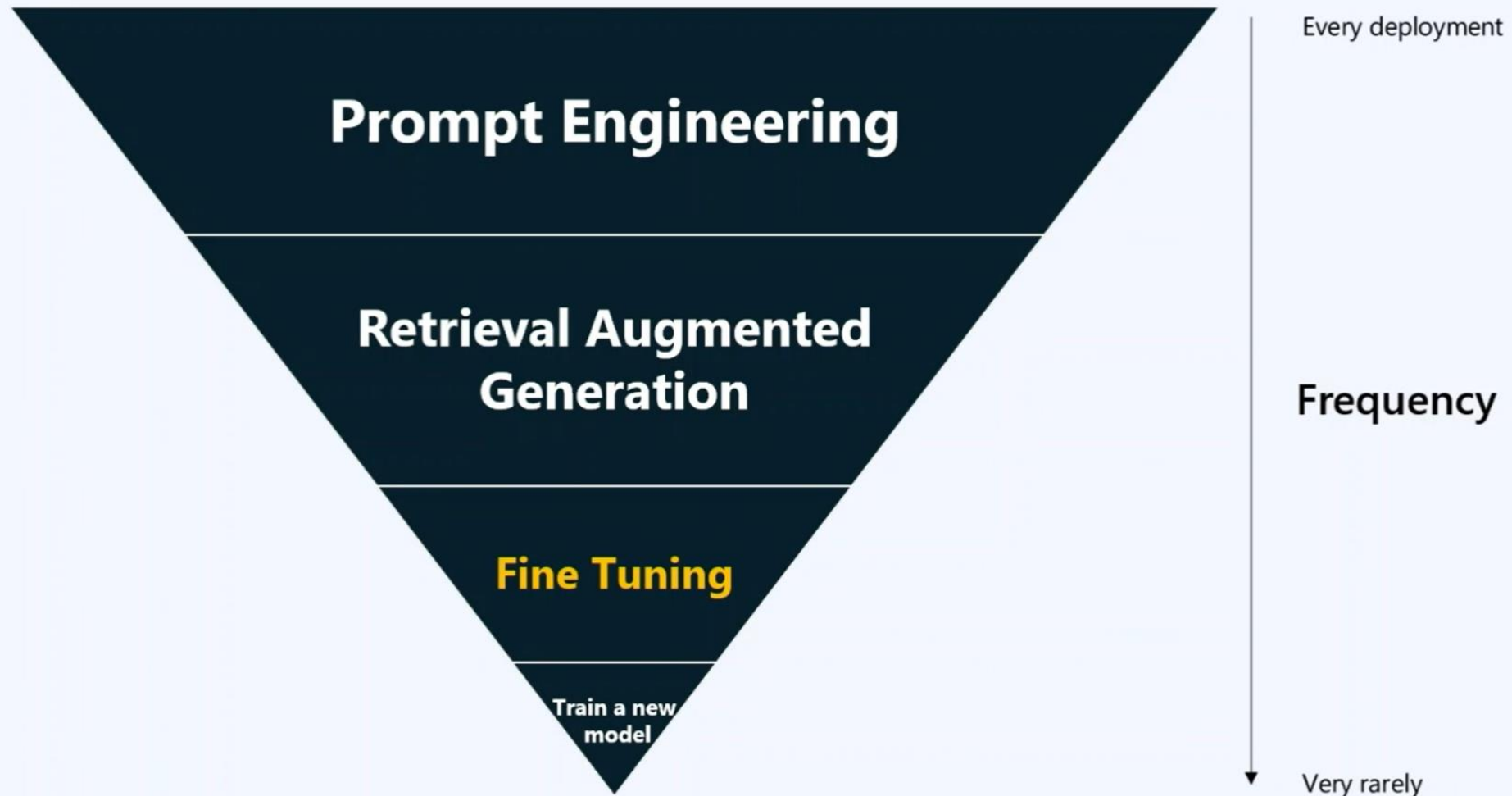Stable max latency and throughput for uniform workloads

**Reserved Processing Capacity**

Ensures capacity is available to customer meet demand.

**Cost Savings**

Potential cost savings for high throughput workloads vs token-based consumption

# Hierarchy of language model customization

# Why fine-tune?

🤩 **Better performance:** Developers hope that by fine tuning models with their own data and instructions, they'll get better results for their tasks

💸 **Cheaper or faster models:** You may want to fine tune a smaller model for a specific task, instead of using an expensive general purpose model like GPT4
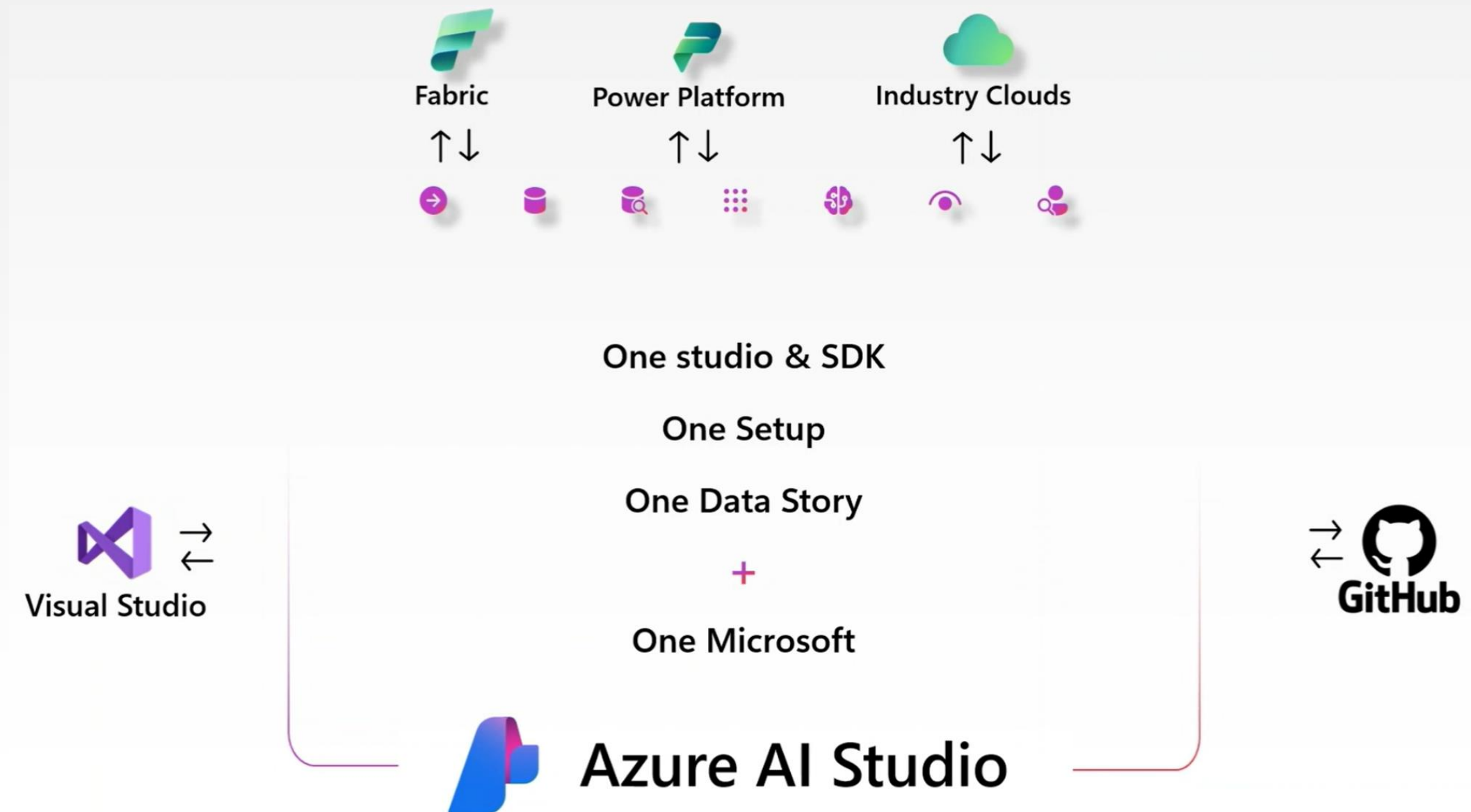
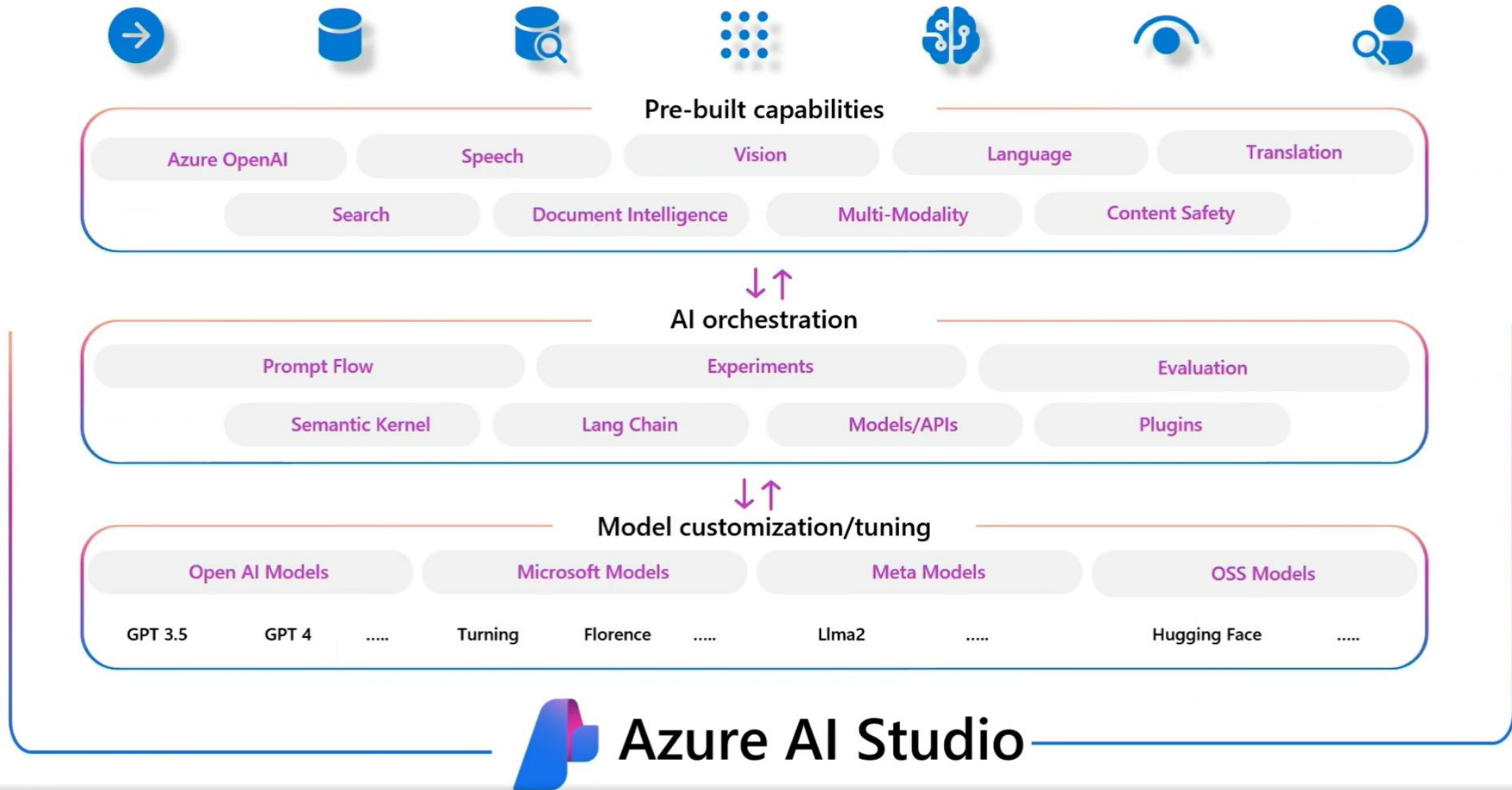❄️ **Differentiation:** Most people won't train a foundation model from scratch; fine tuning with proprietary data provides a competitive advantage

# Model customization cheat sheet

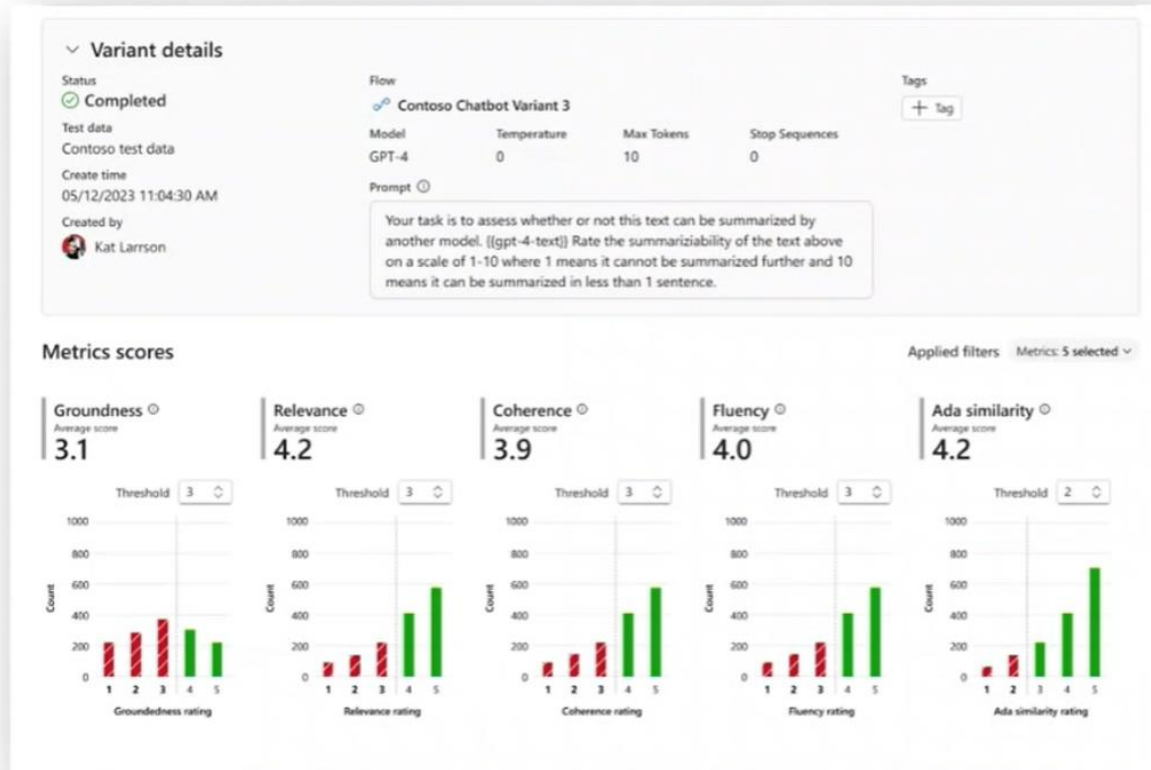| Requirement | Start with | Why? |
|---|---|---|
| **Steer model with a few examples** | Prompt engineering | Easy to craft and quick experimentation, very low barrier to entry |
| **Simple & quick implementation** | Prompt engineering, RAG | Easy tooling with Azure OpenAI on Your Data, PromptFlow, LangChain |
| **Improve model relevancy** | RAG | Retrieve relevant information from your own datasets to insert into prompts |
| **Up to date information** | RAG | Query up to date information from your own databases, search engineers, etc. to insert into prompts |
| **Factual grounding** | RAG | Ability to reference & inspect retrieved data |
| **Optimize for specific tasks** | Fine tuning | Fine tuning is great at steering your model for specific tasks like summarizing data in a specific format |
| **Instructions won't fit in a prompt** | Fine tuning | Fine tuning moves few-shot examples into the training step but increases the quantity of examples are needed to train. |
| **Lower costs** | It depends | ⚠️ Prompt engineering & RAG have lower upfront costs but long prompts are more expensive; training for FT is expensive but may cut prompt length. The choice will always depend on the use case & data. |
| **Complex, novel data or domains** | Prompt Engineering + RAG+ Fine Tuning | ⚠️ This is a high risk area. Fine tuning can retrain the model to recognize new domains, but RAG is needed to avoid plausible confabulations. Make sure customers don't try to retrain for unapproved uses! |

# Azure AI Studio umbrella

# Azure AI Studio umbrella

## Pre-built capabilities

| Azure OpenAI | Speech | Vision | Language | Translation |
| --- | --- | --- | --- | --- |
| | Search | Document Intelligence | Multi-Modality | Content Safety |

↓↑

## AI orchestration

| Prompt Flow | Experiments | Evaluation |
| --- | --- | --- |
| Semantic Kernel | Lang Chain | Models/APIs | Plugins |

↓↑

## Model customization/tuning

| Open AI Models | Microsoft Models | Meta Models | OSS Models |
| --- | --- | --- | --- |
| GPT 3.5    GPT 4    ..... | Turning    Florence    ..... | Llma2    ..... | Hugging Face    ..... |

Azure AI Studio

# Compare model outputs for your business requirements



Evaluate model capabilities, cost, latency, and compatibility with the enterprises' Azure tenant

# Well-architected for reliability & security



Implement a best practice landing zone to ensure you meet the Resilience, Redundancy, and Security needs for your Gen AI implementation

Learn more here:
- [Azure Well-Architected Framework perspective on Azure OpenAI Service - Microsoft Azure Well-Architected Framework](#)
- [Baseline OpenAI end-to-end chat reference architecture - Azure Reference Architectures](#)

# Assistants stack

# BonBon – NashTech intelligent virtual assistant

# BonBon's architecture