

Al at Scale

Introduction to Foundation models in Azure

Keita Onabuta

Senior Customer Engineer for Al & Machine Learning Azure CXP – FastTrack for Azure Microsoft Corporation

Logistics

- Duration about 45min
- Question please post questions in the chat
- This slide deck will be shared after the session
- No recording
 - If you need our support to build and deploy Azure solution, please let us know.

Agenda

- 1. Introduction to Foundation models
- 2. Foundation models in Azure
 - Demo #1 : Prompt Flow
 - Demo #2 : Training large scale model
- 3. Enterprise Search
 - Demo #3: Enterprise Search with Azure Al
- 4. Recap & Call to action

1. Introduction to Foundation models

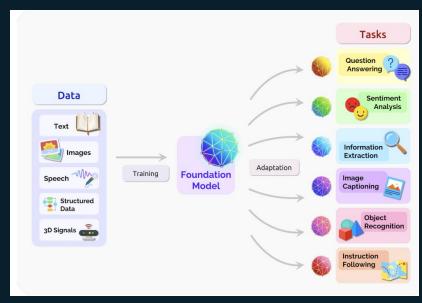
Foundation models

What are Foundation models?

 In recent years, a new successful paradigm for building AI systems has emerged: Train one model on a huge amount of data and adapt it to many applications. We call such a model a foundation model.

Why do we care?

 Foundation models have demonstrated impressive behavior, but can fail unexpectedly, harbor biases, and are poorly understood. Nonetheless, they are being deployed at scale.



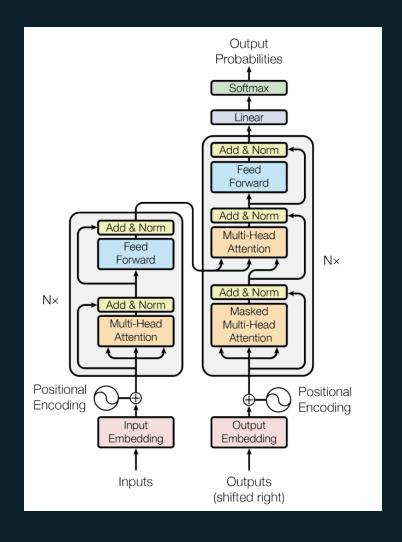
arxiv.org/pdf/2108.07258.pdf

Critical components for Foundation models

- Transformers
- Scale
- · In-context learning

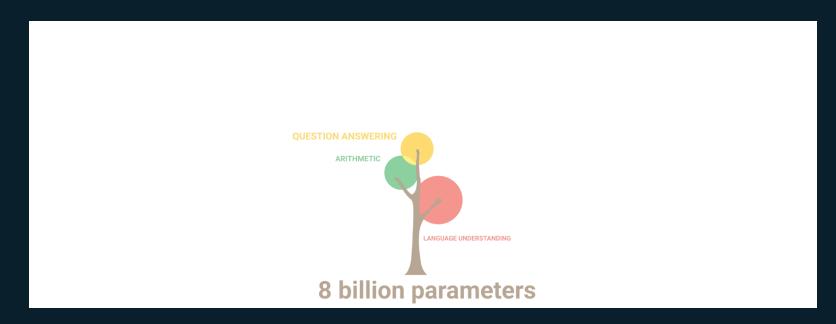
Transformers

- Paper : <u>Attention is All You Need</u>
- Easy to scale and parallelize
 - fast training
 - · more data for training
- Dominating the field of NLP and moving over beyond NLP as well



Scale

- · Scale leads to emerging capabilities
- Many capabilities emerge unpredictably only whe models reach a critical size.



In-context learning

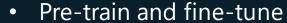
Deep Learning (Representation learning)



• Fully supervised

• architecture design

Pre-trained models (transfer learning)



• no architecture design



Large-scale models (in-context learning)

- Pre-train and prompt
- Zero/Few-shot in-context learning

Why in-context learning?

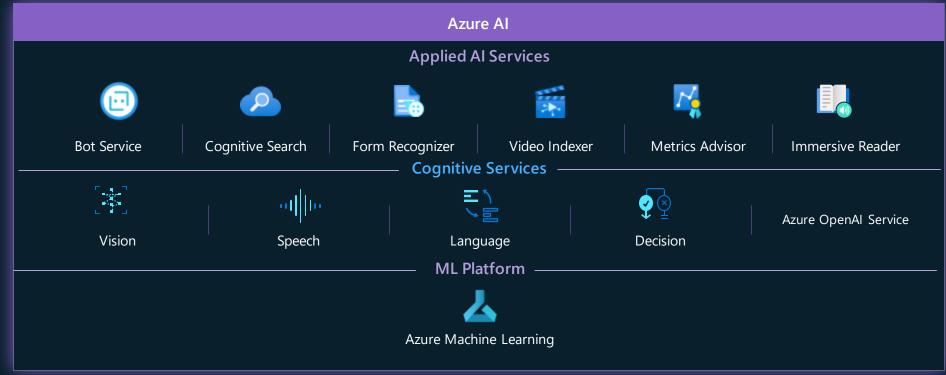
- Models are applied to new tasks out of the box.
- Amazing performance with no or few examples.
- Tasks are adapted to models instead of models adapting to tasks.
- Humans can interact with the models in natural language.
- Blurring the line between ML users and developers.

2. Foundation models in Azure

Microsoft Al portfolio







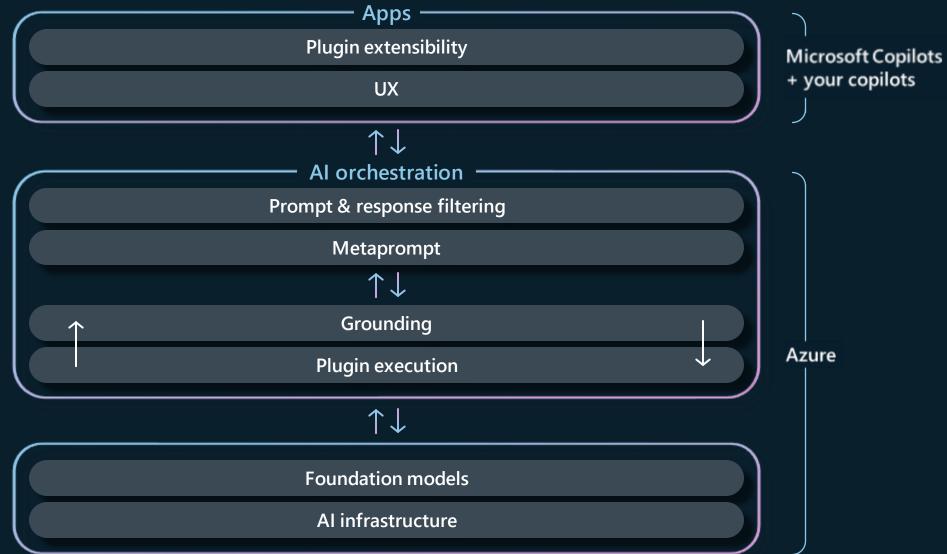


Developers & Data Scientists

Microsoft runs on Azure Al

- · Microsoft 365 Copilot
- Microsoft Security Copilot
- Dynamics 365 Copilot
- Power Platform Copilot
- · Bing chat
- GitHub Copilot
- Nuance
- · LinkedIn
- · And more!

Copilot stack



Microsoft outlines framework for building Al apps and copilots; expands Al plugin ecosystem - Source

Foundation models in Azure Al

Hosted Foundation models

- Azure OpenAl Service
 - Fine tuning supported for some modelss.
- Azure Machine Learning Model Catalog
 - Open source models and Huggin Face models

BYO Foundation models

- Azure Machine Learning
 - Azure Container for PyTorch

Azure OpenAl Service

Large pretrained foundation AI models custom-tunable with your parameters and your data



Summarization Reasoning over data



Writing tools
Code generation



ChatGPT
The Era of Copilots



GPT-3 (GA)

DALL•E 2 (preview)

ChatGPT (GA)

GPT-4 (GA)

Foundation of enterprise security, privacy and compliance

Models

GPT-3

Davinci

- summarizing for specific audience
- Ganerating creative content

Curie

- Answering questions
- Complex, nuanced classification

Babbage

- Semantic search ranking
- Moderately complex classification

Ada

- Simple classification
- Parsing and formatting text

GPT-3.5

GPT 3.5 Turbo

- Primary Chat
- Summarizing
- Code generation

GPT-4

GPT 4 - 32k

• Evaluation of GPT 3.5

GPT 4 - 8k

• Evolution of GPT 3.5

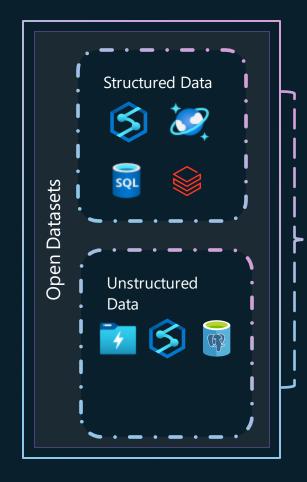
Azure OpenAl | Microsoft Learn

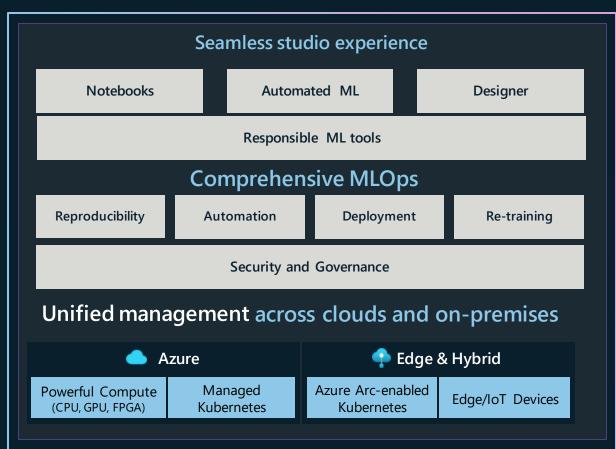
Recent updates (Microsoft Build 2023)

- Azure OpenAl Service on your data (Public Preview)
- Plugins for Azure OpenAl Service (Coming soon)
- Configurable Content Filters
- Provisioned Throughtput (Limited Availablity in June)

Azure Machine Learning

Al Platform for data scientists, machine learning engineers and prompt engineers!







Al-first toolchain

Features for Foundation models

- Compute Cluster
- Model Catalog
 - Open source models and Hugging face models
- Prompt Flow
- Scalable deployment Managed Online Endpoint & Batch Endpoint
- Azure Container for PyTorch
- Model monitoring

Model Catalog (Public Preview)

Hub for Foundation models in AzureML

Open Source Models

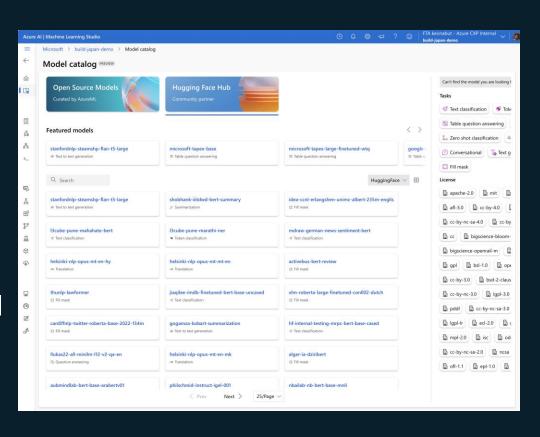
 The most popular open source third-party models curated by Azure Machine Learning. Evaluate, fine tune and deploy models for out of the box usage and are optimized for use in Azure Machine Learning.

Hugging Face hub

 Thousands of models from HuggingFace hub for real time inference with online endpoints.

Azure OpenAI Service

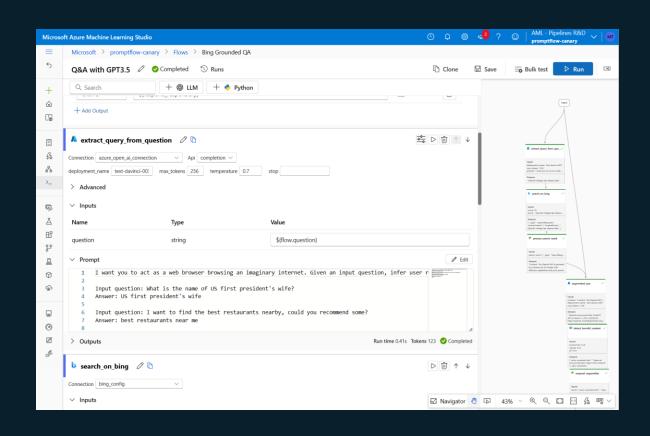
coming soon



Prompt flow (Private Preview)

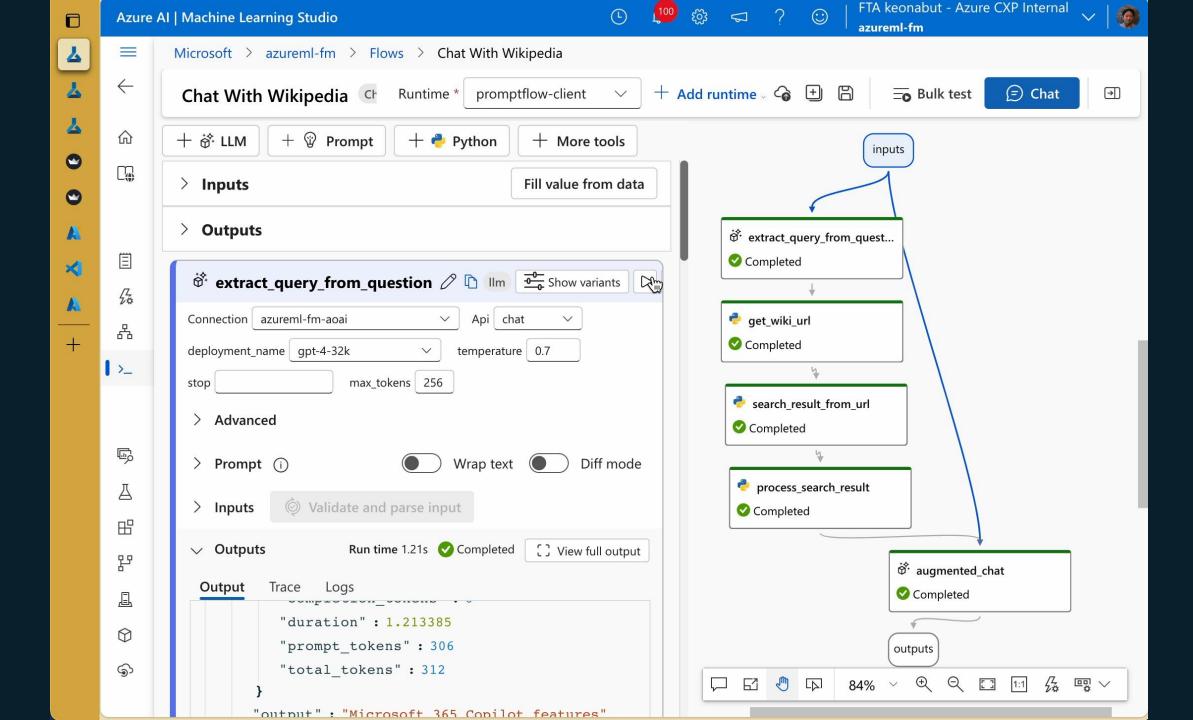
devops for prompt engineering

- Create AI workflows that consume various language models and data sources using the frameworks and APIs of your choice
- One platform to quickly iterate through build, tune, & evaluate for your GenAl workflow
- Evaluate the quality of AI workflows with pre-built and custom metrics
- Easy historical tracking and team collaboration
- Easy deployment and monitoring



Demo #1: Prompt flow

- flow "Chat with Wikipedia"
- Chat with this flow
- Deploy to Managed Online Endpoint
- Call deployed API from Streamlit application



Azure Container for PyToch



Optimized training framework

Set up, develop, and accelerate PyTorch models on large workloads



Up-to-date stack

Latest compatible versions of Ubuntu, Python, PyTorch, Cuda\ROCm, etc.



Ease of use

Installed and validated against dozens of Microsoft workloads to reduce setup costs and accelerate time to value.



Latest training optimization technologies

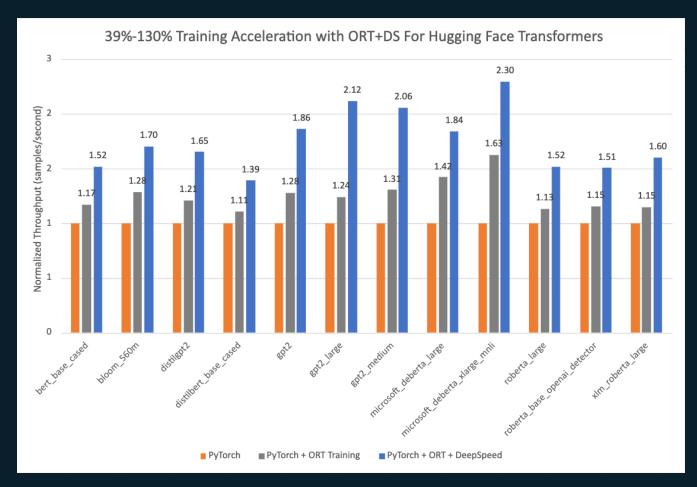
ONNX Runtime, ORT MoE, DeepSpeed, Nebula, MSCCL, and others.



Native integration with Azure

Customer support

Benchmark

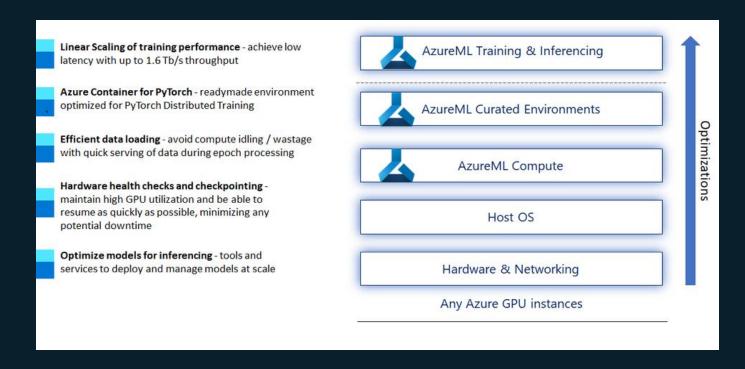


39% - 150% Training Accelereration with ORT+DS For Hugging Face Transformers

Large-scale model in Azure Machine Learning

Best place for high performance deep learning

Microsoft provide best practices for large scale training workloads to get highly efficient optimized performance using state of art technologies.



Key consideration

Azure ML Datastore

- Written in Rust (high speed and high memory efficiency, Avoid issues with Python GIL)
- Multi-process (parallel) data loading etc

Linear scaling with Infiniband Enabled SKUs

• Most of the HPC VM sizes & N-series size designated with 'r' are RDMA-capable.

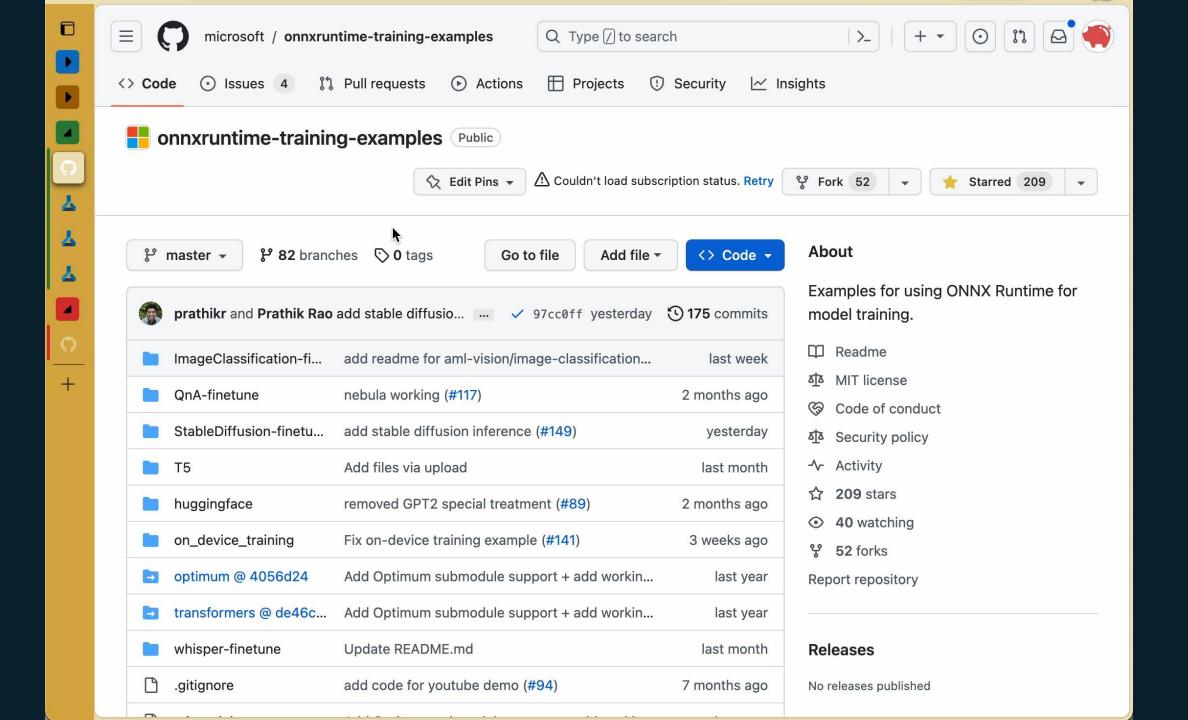
DeepSpeed & ONNXRuntime Training for training optimization

Nebula for fast checkpoint

DeepSpeed-MII for optimized inference

Demo #2: Train large-scale model

- DeepSpeed configuration file
- Python code for finetuning model
- Azure Container for PyTorch (ACPT) as a curated Environment
- Creating custom Environment based on ACPT
- Training the large-scale model with custom Environment



3. Enterprise Search

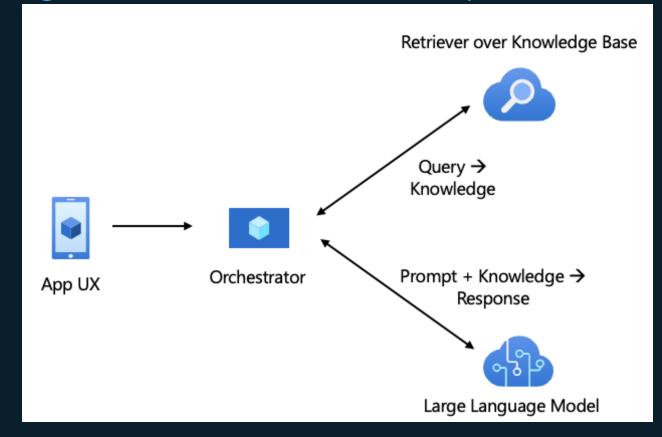
Example of application

- Enterprise search
- Code generation or transformation
- Robotics
- Writing ad
- And more!

Retrieval Augmented Generation

Complement LLMs knowledge by retrieving information relevant to the question

- 1. Find the most relavant infromation from a large data.
- Inline these information in a prompt along with instructions and the question itself.

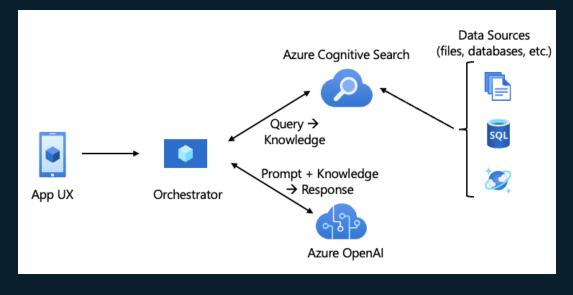


Enterprise search with Azure Cognitive Search

Azure Cognitive Search

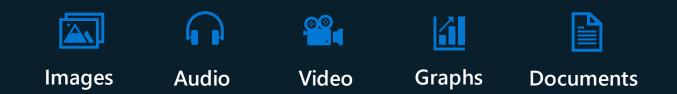
- Azure's complete retrieval solution
- Data ingestion, enterprise-grade security, partitioning and replication for scaling, support for 50+ written languages, and more

Basic architecture



Vector search in Azure Cognitive Search (Private Preview)

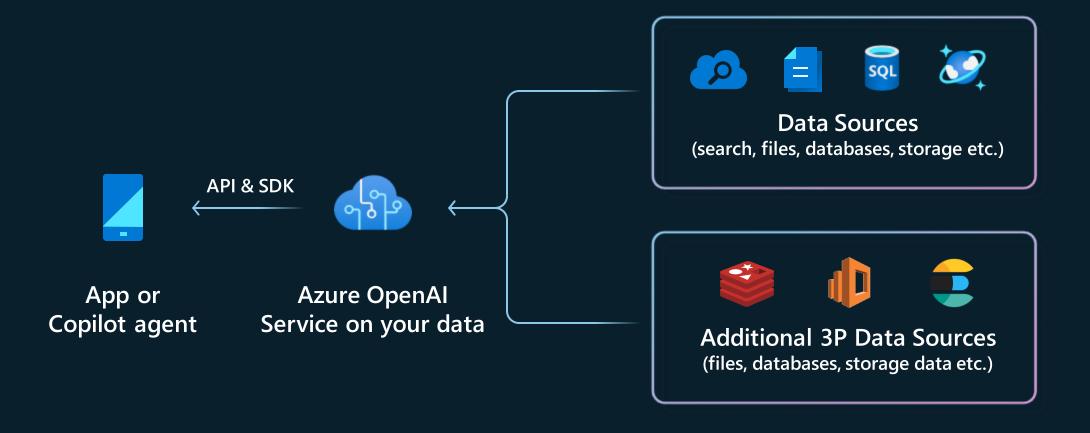
Power your retrieval-augmented generation applications



- Use vector or hybrid search
- Use Azure OpenAI embeddings or bring your own
- Deeply integrate with Azure
- Scale with replication and partitioning
- Build generative Al apps and retrieval plugins

Azure/cognitive-search-vector-pr: The official documentation and code samples for the Vector search feature (preview) in Azure Cognitive Search. (github.com)

Azure OpenAl Service on your data (Public Preview)



Demo #3: Enterprise search

4. Recap & Call to action

Recap

- 1. Introduction to Foundation models
- 2. Foundation models in Azure
 - Demo #1 : Prompt Flow
 - Demo #2 : Training large scale model
- 3. Enterprise Search
 - · Demo #3: Enterprise Search with Azure Al
- 4. Recap & Call to action

Call to action!

Learn Build skills on Microsoft Learn

Develop AI solutions with Azure OpenAI -Training | Microsoft Learn

Connect Join the Al Tech Community to connect, learn, and engage with thousands of members around the world

> Artificial Intelligence and Machine Learning - Microsoft Community Hub

Explore

Stay up to date with the latest news, announcements and release notes

Azure updates | Microsoft Azure

Release note for CLI v2 and Python SDK v2.

Join Use FastTrack for Azure program to accelerate your project!

FastTrack for Azure – Technical Enablement FAQ | Microsoft Azure

Q&A

Thank you!