# Introduction to Azure OpenAI Service

Unlocking the Power of Language Models with Microsoft Azure

Date: April 7, 2025

Ivan Portilla

STAT 5350/4350

# What is Azure OpenAI Service?

- REST API access to OpenAI models via Azure
- Supports GPT-4o, GPT-4 Turbo, GPT-3.5, DALL-E, Whisper, etc.
- Use cases:
  - Text & image generation
  - Summarization
  - Code generation
  - Semantic search
  - Image understanding

https://learn.microsoft.com/en-us/azure/ai-services/openai/overview

# Key Features Overview

- Models: GPT-4o, GPT-4 T, GPT-3.5, Embeddings
- Vision support: GPT-4 Turbo with Vision
- Fine-tuning: GPT-4, GPT-3.5-Turbo (preview)
- UI: Azure Portal, AI Foundry
- Content filtering: Built-in

# Responsible AI with Microsoft

- Microsoft applies Responsible AI principles:
  - Fairness, Safety, Privacy, Transparency
- Tools:
  - Content filters
  - Responsible AI documentation
  - Code of Conduct
  - Limited Access for sensitive features

# How to Get Started

1. Create Azure OpenAI Resource

2. Deploy a Model (e.g., GPT-4o)

3. Use the model:

   – AI Foundry Playground

   – REST API or SDKs (Python, C#, JS, etc.)

# Lab 1 – Set up Az OpenAI service

https://learn.microsoft.com/en-us/azure/ai-services/openai/how-to/create-resource

1. Sign in to Portal

https://portal.azure.com/

1. Create a resource in Azure Services

2. Search Azure OpenAI

# Lab 1 – Set up Az OpenAI service

3. Complete Basics tab (take defaults and click create)

**Create Azure OpenAI** ...

① **Basics** ② Network ③ Tags ④ Review + submit

Azure OpenAI Service provides access to OpenAI's powerful language models, including all the latest OpenAI models. These models can be easily adapted to your specific tasks, including but not limited to content generation, summarization, image understanding, semantic search, and natural language to code translation. Top use cases include Call Centers, Virtual Assistants, Accessibility, Content Generation, and Code Development. The service also features the Assistants API, Fine Tuning capabilities and many ways to connect your data to the service for conversational experiences. The service can be scaled through Standard (tokens) and Provisioned (PTUs) deployment types.

Learn more

**Project Details**

Subscription * ⓘ | Pay-As-You-Go

Resource group * ⓘ | AI-search
Create new

**Instance Details**

Region ⓘ | East US

Name * ⓘ | 5350

Pricing tier * ⓘ | Standard S0

# Lab 2 – Chat playground

- Open AI OpenAI Resource (ai.azure.com)
- In Chat Playground, create Deployment

# Lab 2 – Chat playground

- Select gpt-4o-mini
- Click Confirm
- Accept default values
- Click Deploy

# Lab 2 – Chat playground

- Play with different System prompts, and Top P and Temperature parameters
- Discuss your results

# Lab 3 – Images

- Click on the Images link
- Create Deployment
- Select dall-e-3, click confirm
- Accept Default values
- Click Deploy

# Lab 3 – Images

- Create an image to match sample posted in Slack



## Images playground

| Deployments | dall-e-3 ⌄ | | 🗨 Filters feedback |

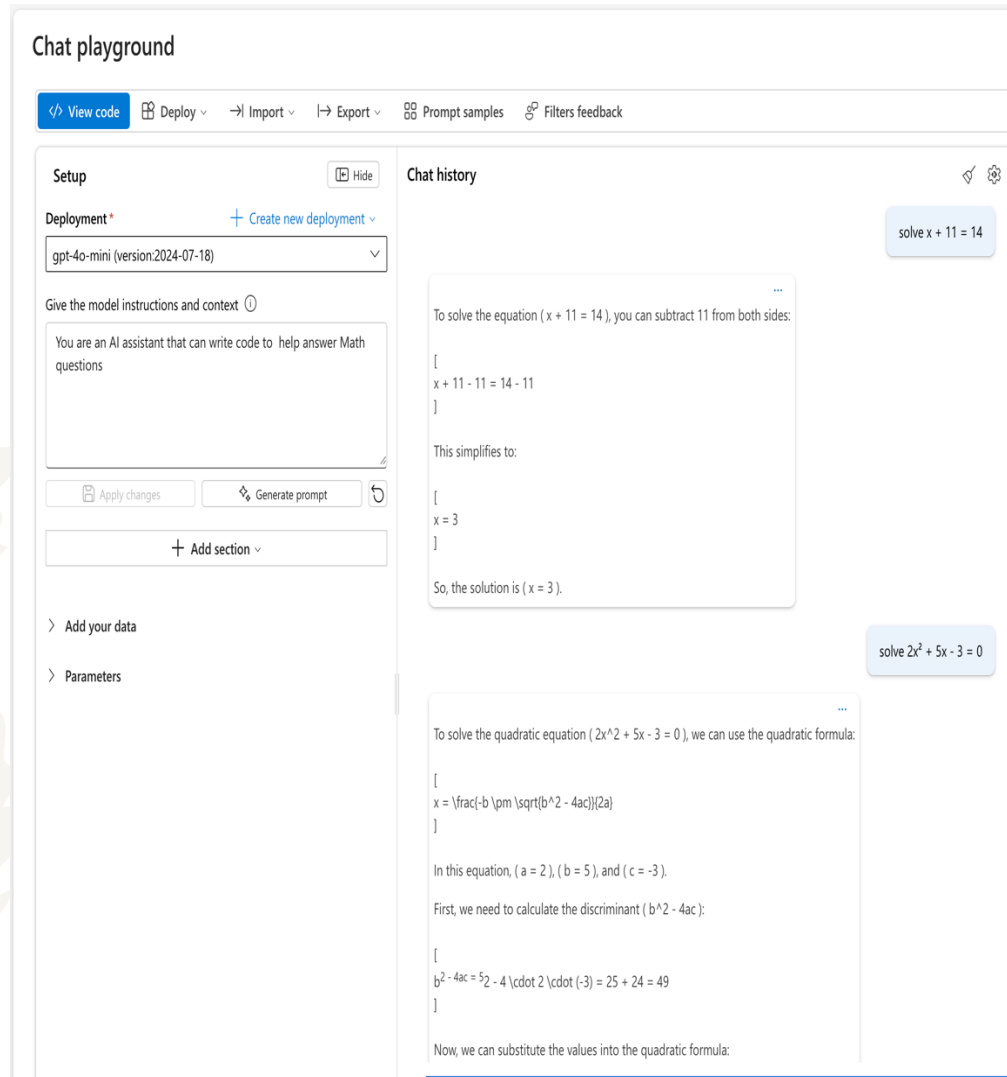**Prompt** ⓘ

lightly sketched image of a bull, Picasso style



lightly sketched image of a bull, Picasso style

# Lab 4 – Assistants

- Same steps
- Solve a quadratic equation with Code Interpreter



https://learn.microsoft.com/en-us/azure/ai-services/openai/assistants-quickstart

# Lab 4 – Assistants