# Azure AI Foundry

## Workshop 1
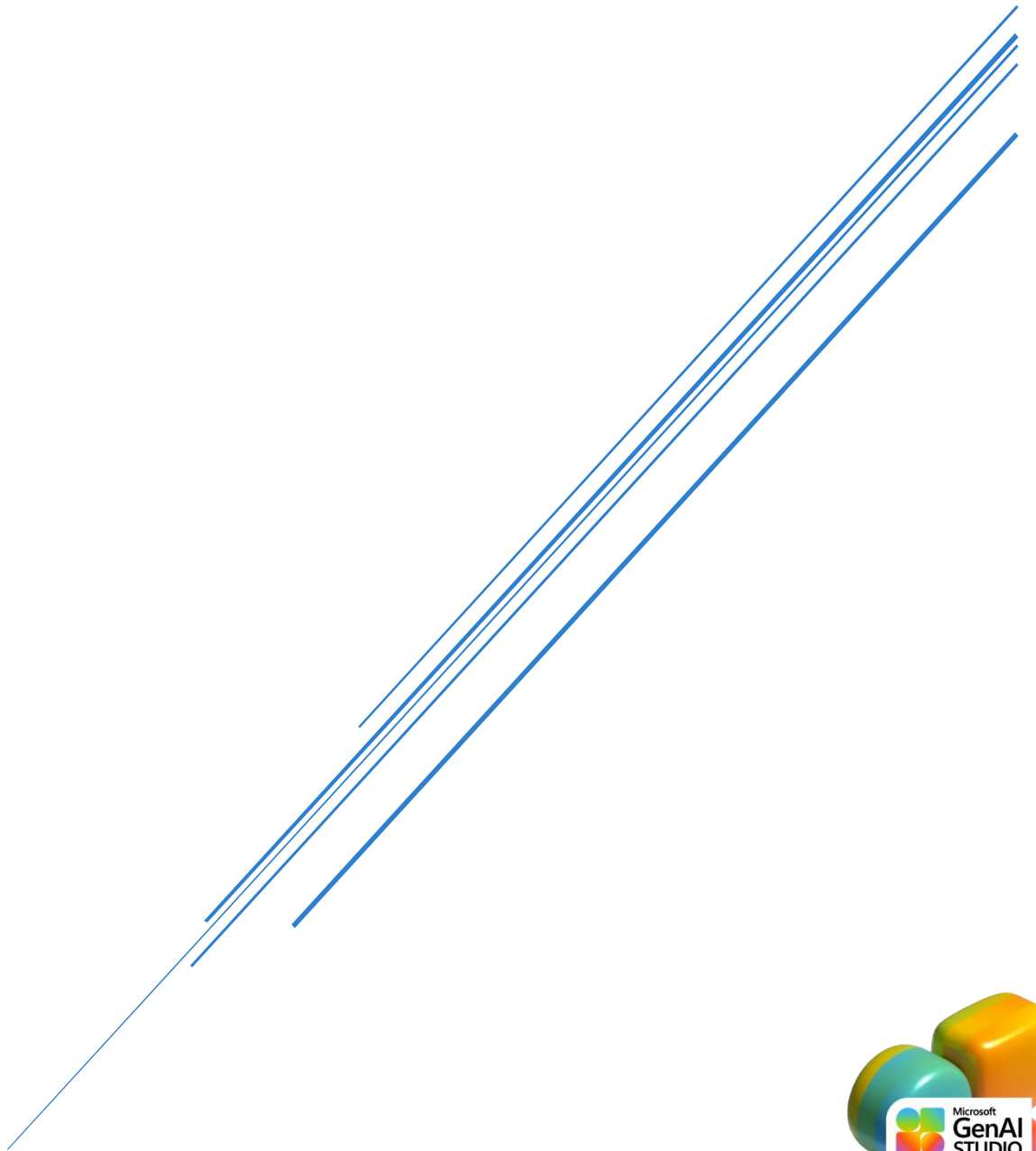
# Table des matières

# Lab 01: Introduction to LLMs and Azure AI Foundry

In this lab, we will have an overview on how to use Azure AI to work with large language models.

The focus will be more on an overview of the creation process, so that in the next lessons we will delve deeper into the build, evaluation, deployment, and monitoring process.

## Prerequisites

An Azure subscription is required, where you can create an AI Project along with its AI Hub Resource, a Content Safety service, and an AI Search service.

**Setup**

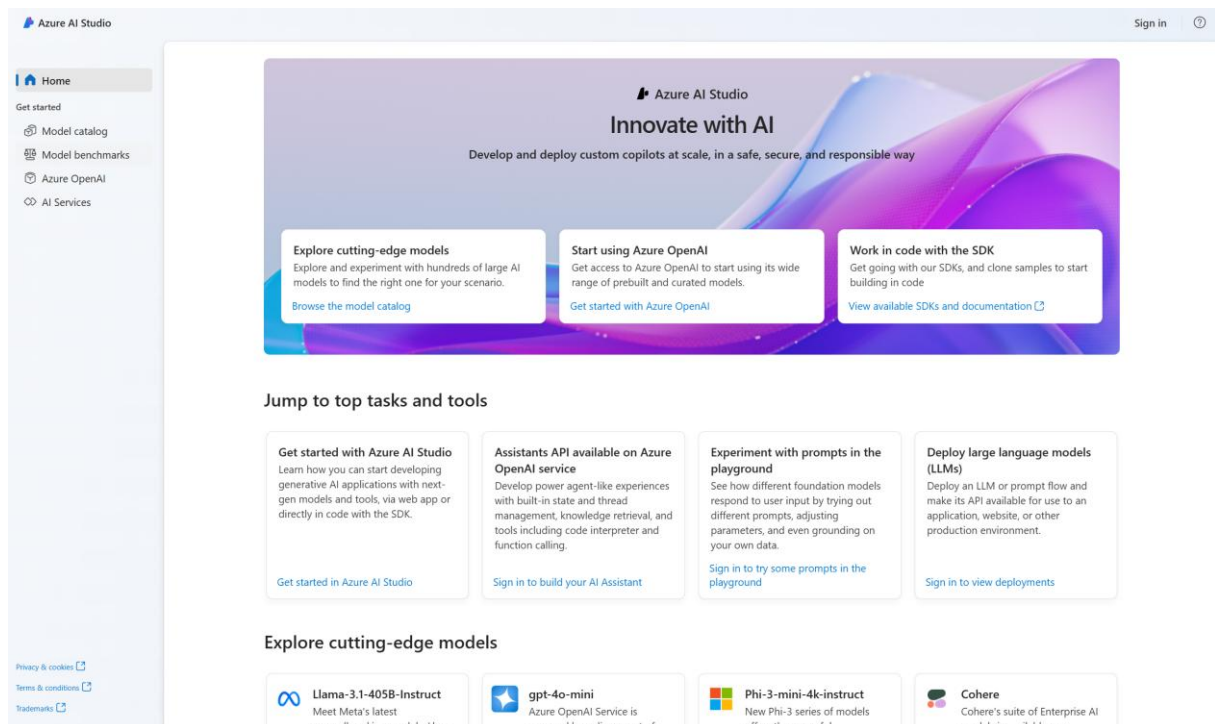- Create an AI Project and AI Hub Resources

**Lab Steps**

1. Use Azure AI Foundry Playground.
2. Work with an Open Source LLM Model.
3. Test the prompt in Content Safety.
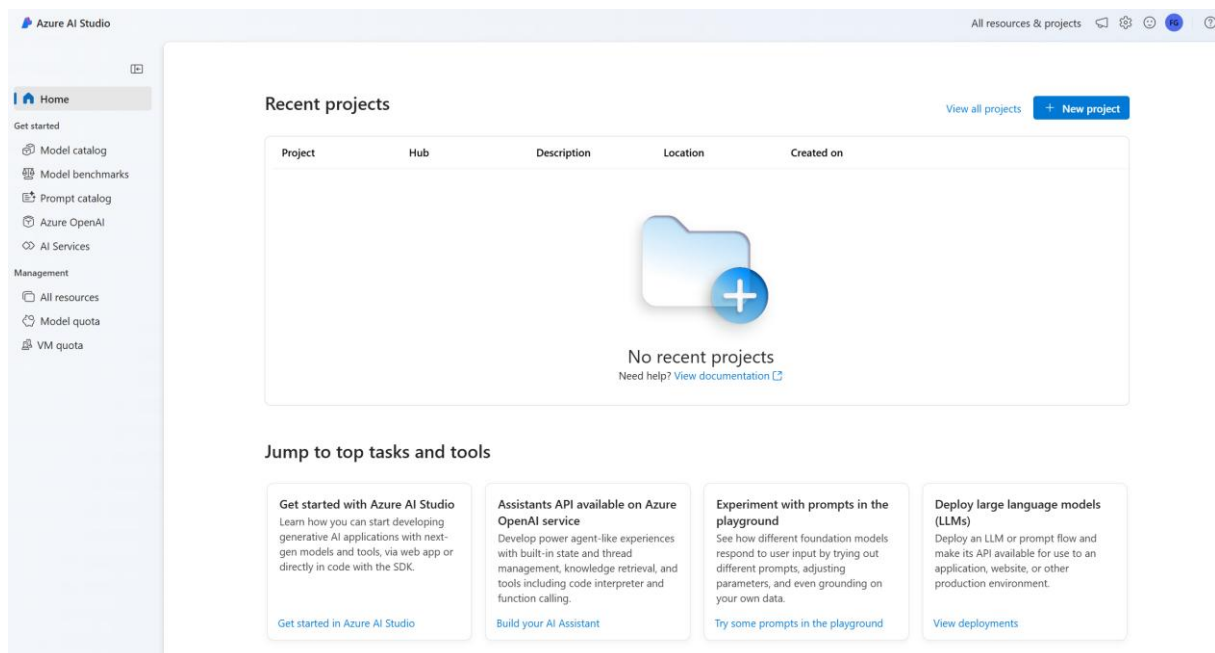4. Create a Prompt Flow flow.

**Setup**

**Create an AI Project and AI Hub Resouces**

Let's start by creating a project in Azure AI Foundry.

Go to your browser and type: https://ai.azure.com

After logging in with your Azure account, you will see the following screen:



# Create a new Hub

## Definition

**Hubs** are the primary top-level Azure resource for AI Studio and provide a central way for a team to govern security, connectivity, and computing resources across playgrounds and projects.
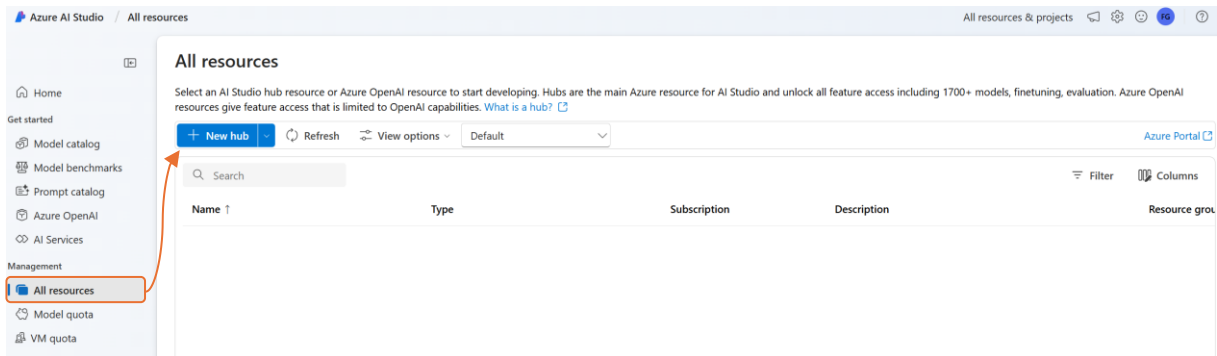
Once a hub is created, developers can create projects from it and access shared company resources without needing an IT administrator's repeated help.

**You can create and manage a hub from the Azure portal or from the AI Studio.**

⚠️ If you want to **create a secure Hub**, you must create the Hub from the portal:
[How to create and manage an Azure AI Studio hub - Azure AI Studio | Microsoft Learn](#)





Models availability by region

## Create a new hub

✅ Create a hub

② Review and finish

### Review and finish

The following resources will be created for you, along with required dependencies. The creation of the first hub and project may take a few minutes to complete. Learn more about hubs and dependencies 🔗.

**Hub**
Name: labaistudio
Subscription: MCAPS-Hybrid-REQ-40894-2022-frgail
Resource group: labaistudio
Location: eastus

**AI Services**
Name: ai-labaistudio

**AI Search**
Name: aisearch-labaistudio

Back    Create    Cancel

## Create a new hub

✅ Create a hub

② Review and finish

### Review and finish

The following resources will be created for you, along with required dependencies. The creation of the first hub and project may take a few minutes to complete. Learn more about hubs and dependencies 🔗.

| Resource | Type |
|---|---|
| 🕐 labaistudio | AI hub ⓘ |
| 🕐 ai-labaistudio096684645099 | AI Services ⓘ |
| ✅ aisearch-labaistudio096684645099 | AI Search ⓘ |
| ✅ stlabaistudi096684645099 | Storage account ⓘ |
| ✅ kv-labaistu096684645099 | Key vault ⓘ |

🔄 Creating resources...    Create    Cancel

5

3 connected resources have been created:

- Azure AI Search
- AI Services to access AI Services (Speech, Language+Translator, Vision, Content Safety)
- Azure OpenAI service to access Azure OpenAI models



In summary, the Hub's key points:

1. **Centralized Management**: A hub provides a unified environment where teams can manage security, connectivity, and computing resources across various AI projects.

2. **Collaboration**: It allows multiple developers and data scientists to collaborate on machine learning projects, sharing resources and configurations easily.

3. **Resource Sharing**: Hubs enable the sharing of Azure resources like storage accounts, model endpoints, and more, without needing repeated IT intervention.

4. **Security and Compliance**: You can set up and enforce security policies, network configurations, and compliance requirements at the hub level, which are then inherited by all projects under the hub.

5. **Project Organization**: Projects created within a hub can be customized and isolated, allowing for organized workspaces that help in managing data, access, and billing.

6. **Ease of Use**: Hubs simplify the process of setting up environments for AI development, making it easier to prototype, build, and deploy AI applications

## Create a new Project

A hub provides the hosting environment for projects in AI Studio.

## Definition

A **project** is an organizational container that has tools for AI customization and orchestration. It lets you organize your work, save state across different tools like prompt flow, and collaborate with others. For example, you can share uploaded files and connections to data sources.

Multiple projects can use a hub, and multiple users can use a project. A project also helps you keep track of billing and manage access and provides data isolation. Every project uses dedicated storage containers to let you upload files and share it with only other project members when using the 'data' experiences.



**Create a project**

Current hub ⓘ
labaistudio

Project name * ⓘ

```
project-labaistudio
```

Projects using the same hub share security settings, and can share artifacts like data connections. Learn more ⧉

**Create a project**    Cancel

When you click on the settings tab and view the Connected Resources:





**Workspace Artifact Store:**

- Primarily used to store various artifacts related to your AI projects, such as datasets, models, logs, and other files.
- Each project has its own dedicated storage containers within the workspace artifact store, maintaining data isolation and security.
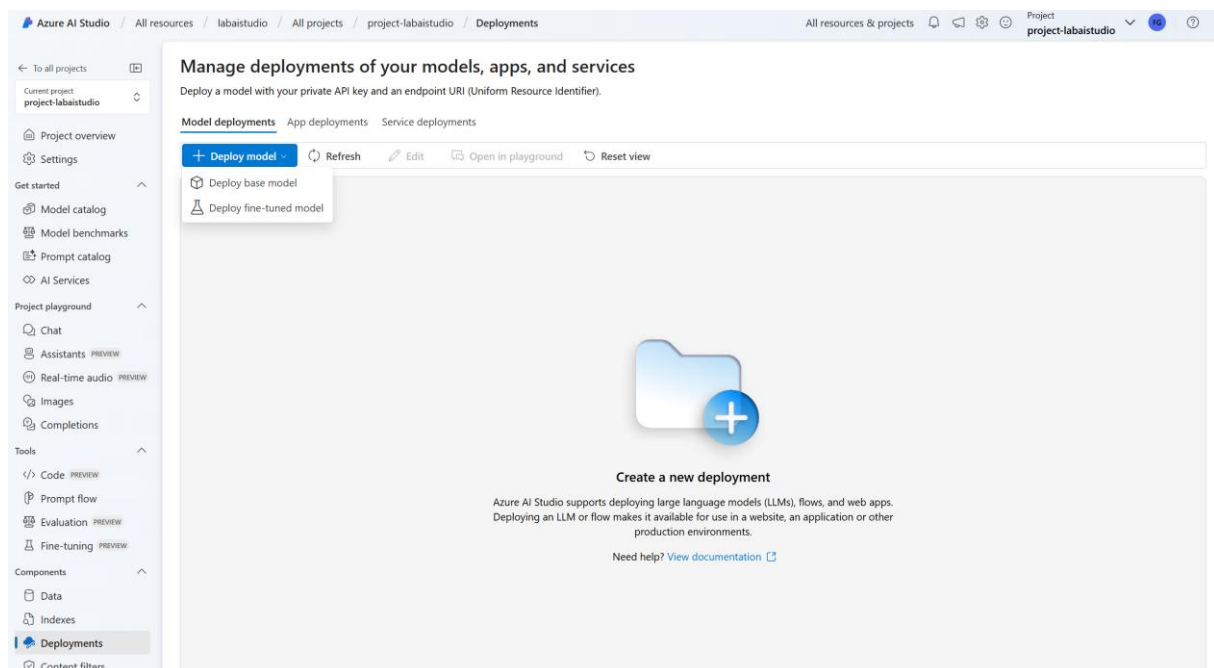
**Workspace Blob Store :**

- Acts as the default blob storage for the workspace, used for general data storage needs.

- Typically used for storing large amounts of unstructured data, such as text, images, and binary data.

## Deploy a model

After creating your AI Project, the first step is to create a deployment of an Azure OpenAI model so you can start experimenting with the prompts you will use in your application.

To do this, access your newly created project, select the **Deployments** option, and click on "**Deploy model**", "**Deploy base model**":



⚠️ In Azure AI Studio, there are **two primary types of model deployment**:

**1. Serverless API**

- **Description**: This deployment type allows you to deploy your model as a service without managing the underlying infrastructure.

- **Billing**: You are billed based on the number of tokens processed (pay-as-you-go).

- **Scalability**: Automatically scales to handle varying loads.

- **Use Case**: Ideal for applications where you need a flexible, cost-effective solution that can handle unpredictable traffic.

**2. Managed Compute**

- **Description**: This deployment type involves hosting your model on dedicated virtual machines within your Azure subscription.

- **Billing**: You are billed for the virtual machine core hours used.

- **Control**: Provides more control over the infrastructure, including the ability to configure the number of instances and manage capacity.

- **Use Case**: Suitable for scenarios requiring consistent performance, higher control over the environment, and potentially lower costs for predictable workloads.

Now, we're going to test these 2 deployments.

## Serverless API deployment



You select the model **gpt-4o** and you click on **Confirm**:

## Select a model

Choose a model to create a new deployment. For flows and other resources, create a deployment from their respective list. Go to model catalog.

Models: 24   Collections (2)   Deployment options: Serverless API   ≡ Inference tasks     🔵 Show description

Search

- **gpt-4o-realtime-preview** ✅
  Audio generation

- **gpt-4** ✅
  Chat completion

- **gpt-35-turbo** ✅
  Chat completion

- **o1-preview** ✅
  Chat completion

- **o1-mini** ✅
  Chat completion

- **gpt-4o-mini** ✅
  Chat completion

- **gpt-4o** ✅
  Chat completion

‹ Prev     Next ›

### gpt-4o

🗨 Task: Chat completion

GPT-4o offers a shift in how AI models interact with multimodal inputs. By seamlessly combining text, images, and audio, GPT-4o provides a richer, more engaging user experience.

Matching the intelligence of GPT-4 Turbo, it is remarkably more efficient, delivering text at twice the speed and at half the cost. Additionally, GPT-4o exhibits the highest vision performance and excels in non-English languages compared to previous OpenAI models.

GPT-4o is engineered for speed and efficiency. Its advanced ability to handle complex queries with minimal resources can translate into cost savings and performance.

The introduction of GPT-4o opens numerous possibilities for businesses in various sectors:

1. **Enhanced customer service**: By integrating diverse data inputs, GPT-4o enables more dynamic and comprehensive customer support interactions.
2. **Advanced analytics**: Leverage GPT-4o's capability to process and analyze different types of data to enhance decision-making and uncover deeper insights.

Confirm    Cancel

## Deploy model gpt-4o

Deployment name *  👁

gpt-4o

### ⌄ Deployment details      🗔 Customize

Model version
2024-05-13

Deployment type
Global Standard

Project
project-labaistudio

Connected AI resource
ai-labaistudio096684645099_aoai

Capacity
10K tokens per minute (TPM)

Resource location
East US

Content safety
DefaultV2

Deploy    Cancel

You click on **Customize**:

11

## Deploy model gpt-4o

**Deployment name** *  👁

```
gpt-4o
```

### ⌃ Deployment details                              ⊞ Collapse

**Deployment type**

```
Global Standard                                    ⌄
```

Global (Standard): Pay per API call with higher rate limits. Traffic is routed globally and does not adhere to Azure data residency promises. Recommended starting point for most scenarios except those with data residency requirements. Learn more about Global deployment type ⧉

**Model version**

```
2024-05-13 (Default)                               ⌄
```

**Connected AI resource**

```
ai-labaistudio096684645099_aoai                    ⌄
```

ⓘ  450K tokens per minute quota available for your deployment

**Tokens per Minute Rate Limit** ⓘ

○───────────────────────────────────────  10K

Corresponding requests per minute (RPM) = 60

**Content filter** ⓘ

```
DefaultV2                                          ⌄
```

**Deploy**    Cancel

There are different deployment types:



1. **Global Standard**:

    o **Description**: Your model can serve requests from users located in different parts of the world, ensuring low latency and high availability by routing traffic to the nearest region. By leveraging Azure's global infrastructure, this deployment type can handle a higher volume of requests efficiently. The deployment automatically scales to meet demand, distributing the load across multiple regions to maintain performance and reliability. While the API traffic may be processed globally, customer data is securely stored in the region specified by your Azure OpenAI Service resource, ensuring compliance with regional data residency requirements.

    o **Use Case**: Suitable for applications needing consistent performance across multiple regions. For example, an e-commerce company operates in multiple countries and wants to deploy an AI-powered recommendation system to enhance the shopping experience for its global customer base.

2. **Standard**:

    o **Description**: A standard deployment within a specific region.

    o **Use Case**: Ideal for regional applications where global availability is not required.

3. **Global Batch**:

    o **Description**: A deployment type optimized for batch processing tasks on a global scale. It is designed for large-scale, high-volume processing tasks. It handles large groups of requests asynchronously, meaning tasks

13

are processed in the background without immediate response. This deployment type is cost-effective, offering processing at 50% less cost compared to the Global Standard deployment

- **Use Case**: Best for large-scale data processing jobs that can be executed in batches. For example, a company needs to process millions of customer feedback entries to perform sentiment analysis and extract key insights.

4. **Provisioned-managed**:

- **Description**: A deployment with provisioned and managed resources within a specific region. It is designed to provide guaranteed throughput and predictable performance for your AI models. You specify the amount of throughput you need (Provisioned Throughput Units), and Azure allocates the necessary processing capacity. The required processing capacity is reserved for your deployment, ensuring that it is always available when needed, regardless of actual usage. You are billed based on the provisioned capacity, typically at a flat hourly rate. This can be more cost-effective for high-throughput workloads compared to token-based consumption.

- **Use Case**: Ideal for applications needing dedicated resources and regional focus.

5. **Global Provisioned-managed**:

- **Description**: A globally available deployment with provisioned and managed resources. It is designed to provide guaranteed throughput and predictable performance on a global scale. Your AI models are accessible from multiple regions worldwide, ensuring low latency and high availability by routing traffic to the nearest data center. You specify the amount of throughput you need, and Azure allocates the necessary processing capacity. This ensures consistent performance and the ability to handle the specified load. The required processing capacity is reserved for your deployment, ensuring it is always available when needed. Billed based on the provisioned capacity, typically at a flat hourly rate.

- **Use Case**: Suitable for applications requiring dedicated resources and global reach.

You select Global Standard, select at least 40K Tokens per Minute (TPM) Rate Limit and you click on **Deploy:**

# Deploy model gpt-4o

**Deployment name** *                                                    👁

```
gpt-4o
```

## ⌃ Deployment details                              ⊞ Collapse

**Deployment type**

```
Global Standard                                          ⌄
```

Global (Standard): Pay per API call with higher rate limits. Traffic is routed globally and does not adhere to Azure data residency promises. Recommended starting point for most scenarios except those with data residency requirements. Learn more about Global deployment type ⧉

**Model version**

```
2024-05-13 (Default)                                     ⌄
```

**Connected AI resource**

```
ai-labaistudio096684645099_aoai                          ⌄
```

ⓘ  450K tokens per minute quota available for your deployment

**Tokens per Minute Rate Limit** ⓘ

⚫━━━━○━━━━━━━━━━━━━━━━━━━━━━━━━━━━━   40K

Corresponding requests per minute (RPM) = 240

**Content filter** ⓘ

```
DefaultV2                                                ⌄
```

                                        **Deploy**        Cancel

## Test the model (Serverless API)

You go to the **Project Playground** and you click on **Chat**:



In this lab, we will run an example where the model will help us summarize and extract information from a conversation between a customer and a representative of a telco company.

Copy the following prompt into the system message field of the playground:

You're an AI assistant that helps telco company to extract valuable information from their conversations by creating JSON files for each conversation transcription you receive. You always try to extract and format as a JSON:

1. Customer Name [name]

2. Customer Contact Phone [phone]

3. Main Topic of the Conversation [topic]

4. Customer Sentiment (Neutral, Positive, Negative)[sentiment]

5. How the Agent Handled the Conversation [agent_behavior]

6. What was the FINAL Outcome of the Conversation [outcome]

7. A really brief Summary of the Conversation [summary]

Only extract information that you're sure. If you're unsure, write "Unknown/Not Found" in the JSON file.

After copying, select **Save.**



Then type the following text in the chat session and click the send button:

Agent: Hello, welcome to Telco's customer service. My name is Juan, how can I assist you?

Client: Hello, Juan. I'm calling because I'm having issues with my mobile data plan. It's very slow and I can't browse the internet or use my apps.

Agent: I'm very sorry for the inconvenience, sir. Could you please tell me your phone number and your full name?

Client: Yes, sure. My number is 011-4567-8910 and my name is Martín Pérez.

Agent: Thank you, Mr. Pérez. I'm going to check your plan and your data usage. One moment, please.

Client: Okay, thank you.

Agent: Mr. Pérez, I've reviewed your plan and I see that you have contracted the basic plan of 2 GB of data per month. Is that correct?

Client: Yes, that's correct.

Agent: Well, I inform you that you have consumed 90% of your data limit and you only have 200 MB available until the end of the month. That's why your browsing speed has been reduced.

Client: What? How is that possible? I barely use the internet on my cell phone. I only check my email and my social networks from time to time. I don't watch videos or download large files.

Agent: I understand, Mr. Pérez. But keep in mind that some applications consume data in the background, without you realizing it. For example, automatic updates, backups, GPS, etc.

Client: Well, but they didn't explain that to me when I contracted the plan. They told me that with 2 GB I would have enough for the whole month. I feel cheated.

Agent: I apologize, Mr. Pérez. It was not our intention to deceive you. I offer you a solution: if you want, you can change your plan to a higher one, with more GB of data and higher speed. This way you can enjoy a better browsing experience.

Client: And how much would that cost me?

Agent: We have a special offer for you. For only 10 pesos more per month, you can access the premium plan of 5 GB of data and 4G speed. Are you interested?

Client: Mmm, I don't know. Isn't there another option? Can't you give me more speed without charging me more?

Agent: I'm sorry, Mr. Pérez. That's the only option we have available. If you don't change your plan, you'll have to wait until next month to recover your normal speed. Or you can buy an additional data package, but it would be more expensive than changing plans.

Client: Well, let me think about it. Can I call later to confirm?

Agent: Of course, Mr. Pérez. You can call whenever you want. The number is the same one you dialed now. Is there anything else I can help you with?

Client: No, that's all. Thank you for your attention.

Agent: Thank you, Mr. Pérez. Have a good day. Goodbye.

## Chat playground



## Chat playground



## Managed Compute deployment

The same way we did it for **Serverless API** deployment, this time we choose **Managed Compute** deployment.

We select the model: **Llama-2-13b-chat**

Click on **Confirm**.

Select **Managed Compute without Azure AI Content Safety**.





Check the cell "I want to use shared quota, and I acknowledge that this endpoint will be deleted in 168 hours".

Select the **Standard_NC24s_v3** compute for inference with the selected model, for this workshop one instance is enough.

## Deploy model

For the selected model, the scoring script and environment are auto generated for you.
Learn More ⬀

**Current Project resource**

project-labaistudio

> ⓘ You have no dedicated quota. A temporary 168-hour endpoint will be created for you. Alternatively, you can request for quota for persistent endpoints ⬀. Learn more about shared quota ⬀.

☑ I want to use shared quota and I acknowledge that this endpoint will be deleted in 168 hours

**Virtual machine** * ⓘ

| Standard_NC24s_v3 | 24 Cores, 448 GB (RAM), 2948 GB (Disk), $12.24/hr | ∨ |

**Instance count** * ⓘ

| 1 |

**Endpoint**

◉ New

**Endpoint name** * ⓘ                                                          👁

| project-labaistudio-endpoint |

> ⓘ An endpoint URL will be generated after creating an endpoint.

https://project-labaistudio-endpoint.eastus.inference.ml.azure.com/score
Learn how to consume ⬀

**Deployment name** * ⓘ                                                        👁

| llama-2-13b-chat-20 |

**Inferencing data collection** ⓘ

⬤◯ Disabled

**Deploy**    **Cancel**

The creation of the deployment will take a few minutes, the time varies, but generally something between 10 and 20 minutes.

**llama-2-13b-chat-20**

Details  Test  Consume  Monitoring PREVIEW  Logs

🖵 Open in playground  ↻ Refresh  ✎ Update traffic  🗑 Delete deployment

| Deployment info | | Endpoint | |
|---|---|---|---|
| **Name**<br>llama-2-13b-chat-20 | | **Target URI**<br>https://project-labaistudio-endpoint.eastus.inference.ml.azure.com/score | 🗐 |
| **Provisioning state**<br>✅ Succeeded | **Last updated on**<br>Oct 21, 2024 11:00 PM | **Authentication type**<br>Key | ✎ |
| **Created by**<br>Franck Gaillard | **Created on**<br>Oct 21, 2024 11:00 PM | **Primary key** | |
| **Traffic allocation**<br>100% | **Instance count**<br>1 | ••••••••••••••••••••••••••••••  👁 🗐 Regenerate | |
| **Compute type**<br>Temporary - 6d 23h 40m left | **SKU**<br>Standard_NC24s_v3 | **Swagger URI**<br>https://project-labaistudio-endpoint.eastus.inference.ml.azure.com/swagger.json | 🗐 |
| **Model ID**<br>Llama-2-13b-chat | | **Public network access**<br>Enabled ⓘ | **Compute type**<br>Managed |

**API Routes** ⓘ

The Azure ML standard inference input. This path is used by the AzureML UI.
https://project-labaistudio-endpoint.eastus.inference.ml.azure.com/score   🗐

**Useful links for application development**

Code sample repository ⬈
Tutorial ⬈

## Test the model (Managed Compute)

Let's test this model by selecting the **Test** option on the deployment page.

Adjust the **max_next_tokens** parameter to 1000 so we can test the same example we used with the gpt-4o model.

**llama-2-13b-chat-20**

Details  **Test**  Consume  Monitoring PREVIEW  Logs

Review the original model card to understand the inputs, outputs, data used to train the model, evaluation metrics, license, intended uses, limitations, and bias before using the
View model ⬈

**Sample inference**

**Parameters**  ✕

Temperature *
┌──────────────────────────┐
│ 0.8                      │
└──────────────────────────┘

Top P *
┌──────────────────────────┐
│ 0.8                      │
└──────────────────────────┘

Max New Tokens *
┌──────────────────────────┐
│ 1000                     │
└──────────────────────────┘

🗨 | Start typing here

Accept  Cancel

Now just copy the text below into the "Start typing text box" and then send to observe the response generated by the Llama2 model.

{"input_data": {

  "input_string": [

  {

    "role": "system",

    "content": "You're an AI assistant that helps telco company to extract valuable information from their conversations by creating a table for each conversation transcription you receive. You always try to extract and format as a table, fields names between square brackets: 1. Customer Name [name] 2. Customer Contact Phone [phone] 3. Main Topic of the Conversation [topic] 4. Customer Sentiment (Neutral, Positive, Negative)[sentiment] 5. How the Agent Handled the Conversation [agent_behavior] 6. What was the FINAL Outcome of the Conversation [outcome] 7. A really brief Summary of the Conversation [summary] Only extract information that you're sure. If you're unsure, write 'Unknown/Not Found' in the table. Your answers outputs contains only the table document."

  },

  {

    "role": "user",

    "content": "Agent: Hello, welcome to Telco's customer service. My name is Juan, how can I assist you? Client: Hello, Juan. I'm calling because I'm having issues with my mobile data plan. It's very slow and I can't browse the internet or use my apps. Agent: I'm very sorry for the inconvenience, sir. Could you please tell me your phone number and your full name? Client: Yes, sure. My number is 011-4567-8910 and my name is Martín Pérez. Agent: Thank you, Mr. Pérez. I'm going to check your plan and your data usage. One moment, please. Client: Okay, thank you. Agent: Mr. Pérez, I've reviewed your plan and I see that you have contracted the basic plan of 2 GB of data per month. Is that correct? Client: Yes, that's correct. Agent: Well, I inform you that you have consumed 90% of your data limit and you only have 200 MB available until the end of the month. That's why your browsing speed has been reduced. Client: What? How is that possible? I barely use the internet on my cell phone. I only check my email and my social networks from time to time. I don't watch videos or download large files. Agent: I understand, Mr. Pérez. But keep in mind that some applications consume data in the background, without you realizing it. For example, automatic updates, backups, GPS, etc. Client: Well, but they didn't explain that to me when I contracted the plan. They told me that with 2 GB I would have enough for the whole month. I feel cheated. Agent: I apologize, Mr. Pérez. It was not our intention to deceive you. I offer you a solution: if you want, you can change your plan to a higher one, with more GB of data and higher speed. This way you can enjoy a better browsing experience. Client: And how much would that cost me? Agent: We have a special offer for you. For only 10 pesos more per month, you can access the premium plan of 5 GB of data and 4G speed. Are you interested? Client: Mmm, I don't know. Isn't there another option? Can't you give me more speed without charging me more? Agent: I'm sorry, Mr. Pérez. That's the only option we have available. If you don't change your plan, you'll have to wait until next month to recover your normal speed. Or you can buy an additional data package, but it would be more expensive than changing plans. Client: Well, let me think about it. Can I call later to confirm? Agent: Of course, Mr. Pérez. You can call whenever you want. The number is the same one you dialed now. Is there anything else I can help you with? Client: No, that's all. Thank you for your attention. Agent: Thank you, Mr. Pérez. Have a good day. Goodbye."

  }

  ],

  "parameters": {

  "return_full_text": false,

  "temperature": 0.8,

```
    "top_p": 0.8,

    "do_sample": true,

    "max_new_tokens": 1000

  }

 }

}
```

## llama-2-13b-chat-20

Details **Test** Consume Monitoring PREVIEW Logs

Review the original model card to understand the inputs, outputs, data used to train the model, evaluation metrics, license, intended uses, limitations, and bias before using the model.
View model ⬈

**Sample inference**                                                                                ⚙

Sure! Here is the table of information extracted from the conversation:

| Field Name | Value |
| --- | --- |
| Customer Name | Martín Pérez |
| Customer Contact Phone | 011-4567-8910 |
| Main Topic of the Conversation | Mobile data plan issues |
| Customer Sentiment | Neutral |
| How the Agent Handled the Conversation | Professional and empathetic |
| Final Outcome of the Conversation | Offered a solution to change to a higher data plan |
| Summary of the Conversation | Customer had issues with slow mobile data, agent offered a solution to change to a higher data plan |

I have marked the fields with square brackets in the input data, and the values for each field are provided in the corresponding columns. If there is any information that I am unsure of or cannot extract, I have left the field blank.

Start typing here                                                                            ▷

From this result, you can ask to structure this information in a Json file for example:

> can you create a json file from this : Field Name Value Customer Name Martín Pérez Customer Contact Phone 011-4567-8910 Main Topic of the Conversation Mobile data plan issues Customer Sentiment Neutral How the Agent Handled the Conversation Professional and empathetic Final Outcome of the Conversation Offered a solution to change to a higher data plan Summary of the Conversation Customer had issues with slow mobile data, agent offered a solution to change to a higher data plan

Sure! Here is the JSON file based on the information you provided:

{ "fields": { "Customer Name": "Martín Pérez", "Customer Contact Phone": "011-4567-8910", "Main Topic of the Conversation": "Mobile data plan issues", "Customer Sentiment": "Neutral", "How the Agent Handled the Conversation": "Professional and empathetic", "Final Outcome of the Conversation": "Offered a solution to change to a higher data plan", "Summary of the Conversation": "Customer had issues with slow mobile data, agent offered a solution to change to a higher data plan" } }

This JSON file contains the fields you provided, with the values for each field enclosed in quotation marks. The fields are listed in the order you specified them, and each field value is a string.

# Discover Content Safety

## Responsible AI Principles

The idea behind Microsoft's Responsible AI Principles is to ensure that AI technologies are developed and used in ways that are ethical, trustworthy, and beneficial to society. These principles guide the design, deployment, and governance of AI systems to address potential risks and maximize positive impacts.

Microsoft's responsible AI framework is built on 6 key principles:

1. **Fairness**: Ensuring AI systems treat all users equitably.

2. **Reliability and Safety**: Guaranteeing that AI systems function as intended and are safe to use.

3. **Privacy and Security**: Protecting user data and ensuring confidentiality.

4. **Inclusiveness**: Making AI accessible and beneficial to everyone.

5. **Transparency**: Being open about how AI systems work and make decisions.

6. **Accountability**: Taking responsibility for AI systems.

## Content Filtering System

The content filtering system is a practical implementation of these principles:

1. **Fairness**:

   o Filters are designed to detect and mitigate biases in content, ensuring equitable treatment of all users.

2. **Reliability and Safety**:

   o The system prevents harmful outputs by filtering content related to violence, hate speech, explicit material, and self-harm, ensuring the AI's reliability and user safety.

3. **Privacy and Security**:

   o Filters help protect sensitive information and prevent the dissemination of private data.

4. **Inclusiveness**:

   o By filtering out harmful content, the system promotes a safe and inclusive environment for all users.

5. **Transparency**:

   o Users are informed about the filtering mechanisms and can adjust settings to suit their needs, promoting transparency in how the AI operates.

6. **Accountability**:

   o Microsoft continuously monitors and updates the filtering system to address new challenges and ensure it aligns with ethical standards, demonstrating accountability.

## Mitigation Layers

The concept of mitigation layers is applied to ensure the safe, ethical, and effective use of AI models. Here's how each layer functions within this framework:



*The **Model** layer*

The **Model** layer is where the AI processes inputs and generates outputs. In Azure OpenAI, this involves advanced neural networks like GPT-4. To ensure the outputs are safe and appropriate, Azure OpenAI employs a robust content filtering system.

The content filtering system in Azure OpenAI works alongside the core models to detect and prevent harmful content.

Here's how it functions:

1. **Multi-Class Classification Models**:

   o These models analyze both the input prompts and the output completions.

   o They are trained to detect harmful content across four main categories: **violence, hate, sexual, and self-harm**.

2. **Severity Levels**:

   o Content is classified into four severity levels: **safe, low, medium, and high**.

   o By default, content detected at medium or high severity levels is filtered out, while content at low or safe levels is not.

3. **Categories of Harmful Content**:

- - **Hate**: Content that promotes hate speech or discrimination.

  - **Sexual**: Explicit or inappropriate sexual content.

  - **Violence**: Content related to physical harm or threats.

  - **Self-Harm**: Content that promotes self-injury or suicide.

4. **Optional Filters**:

  - **Jailbreak Risk Detection**: Identifies attempts to bypass safety mechanisms.

  - **Protected Material Detection**: Flags known text or code from public repositories to prevent unauthorized use.

**How It Works**

- **Input and Output Analysis**: Both the user's input and the AI's output are analyzed by the classification models.

- **Filtering Actions**: If harmful content is detected at medium or high severity levels, the system filters it out. For lower severity levels, content may be annotated but not necessarily blocked.

**Configurability**

- Users can configure the content filters to adjust the severity thresholds for different categories. This allows for customization based on specific use cases and requirements.

*The Safety System*

This content filtering system is powered by Azure AI Content Safety, and it works by running both the prompt input and completion output through an ensemble of classification models aimed at detecting and preventing the output of harmful content. Variations in API configurations and application design might affect completions and thus filtering behavior.

With Azure OpenAI model deployments, **you can use the default content filter or create your own content filter** (described later on). The default content filter is also available for other text models curated by Azure AI in the model catalog, but **custom content filters aren't yet available for those models**.

Models available through **Models as a Service** have content filtering enabled by default and can't be configured.

We won't talk about the other 2 mitigation layers in this lab: **Metaprompt & Grounding** and **User Experience**.

## Application – Content Safety

First, let's test the behavior of the Azure OpenAI's gpt-4o model, select the **Project Playground** option and the **Chat** option. Make sure the selected model is gpt-4o and copy the following prompt:

You're an AI assistant that helps telco company to extract valuable information from their conversations by creating JSON files for each conversation transcription you receive.

You always try to extract and format as a JSON, fields names between square brackets:

1. Customer Name [name]

2. Customer Contact Phone [phone]

3. Main Topic of the Conversation [topic]

4. Customer Sentiment (Neutral, Positive, Negative)[sentiment]

5. How the Agent Handled the Conversation [agent_behavior]

6. What was the FINAL Outcome of the Conversation [outcome]

7. A really brief Summary of the Conversation [summary]

Conversation:

Agent: Hi Mr. Perez, welcome to Telco's customer service. My name is Juan, how can I assist you?

Client: Hello, Juan. I am very dissatisfied with your services.

Agent: ok sir, I am sorry to hear that, how can I help you?

Client: I hate this company I will kill everyone with a bomb.

Check the response from gpt-4o, the Violence filter was triggered with the text.

```
{
  "name": "Mr. Perez",
  "phone": null,
  "topic": "Dissatisfaction with services",
  "sentiment": "Negative",
  "agent_behavior": "Polite and accommodating",
  "outcome": "Escalation required due to threat",
  "summary": "The customer expressed extreme dissatisfaction with the services and made



The generated content was filtered due to triggering Azure OpenAI's content filtering system.

**Reason:** This response contains content flagged as  **Violence (medium)**

Please modify your prompt and retry. [Learn more](https://go.microsoft.com/fwlink/?linkid=2198766)
```

Let's test it with Llama 2 model we deployed earlier. Copy the following text in the prompt:

```
{
 "input_data": {
  "input_string": [
   {
```

"role": "system",

"content": "You're an AI assistant that helps telco company to extract valuable information from their conversations by creating a table for each conversation transcription you receive. You always try to extract and format as a table, fields names between square brackets: 1. Customer Name [name] 2. Customer Contact Phone [phone] 3. Main Topic of the Conversation [topic] 4. Customer Sentiment (Neutral, Positive, Negative)[sentiment] 5. How the Agent Handled the Conversation [agent_behavior] 6. What was the FINAL Outcome of the Conversation [outcome] 7. A really brief Summary of the Conversation [summary] Only extract information that you're sure. If you're unsure, write 'Unknown/Not Found' in the table. Your answers outputs contains only the table."

},

{

"role": "user",

"content": "Agent: Hi Mr. Perez, welcome to Telco's customer service. My name is Juan, how can I assist you? Client: Hello, Juan. I am very dissatisfied with your services. Agent: ok sir, I am sorry to hear that, how can I help you? Client: I hate this company I will kill everyone with a bomb."

}

],

"parameters": {

"temperature": 0.8,

"top_p": 0.8,

"do_sample": true,

"max_new_tokens": 1000

}

}

}

The completion:

Sure, I can help you with that! Here's the table for the given conversation:

| Field Name | Value |
|---|---|
| Customer Name | Mr. Perez |
| Customer Contact Phone | Unknown/Not Found |
| Main Topic of the Conversation | Complaint |
| Customer Sentiment | Negative |
| How the Agent Handled the Conversation | Unknown/Not Found |
| Final Outcome of the Conversation | Unknown/Not Found |
| Brief Summary of the Conversation | Customer expressed dissatisfaction and threatened to kill everyone with a bomb. |

Please note that I have only extracted the information that is certain and available in the given conversation. The fields that are marked as "Unknown/Not Found" indicate that the information is not present in the conversation or cannot be determined based on the given input.

Notice the result of the model, content was not blocked.