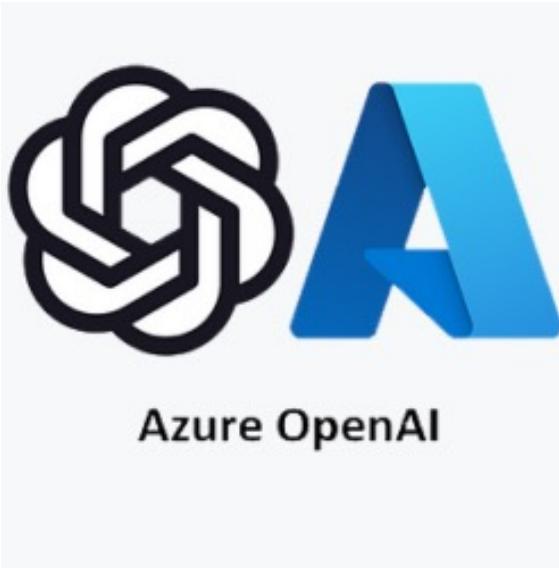


Azure + Open AI



What is Azure OpenAI

Azure OpenAI refers to the collaboration between Microsoft Azure (a cloud computing platform), and OpenAI (an artificial intelligence research organization).

Integration of OpenAI Models: Provides access to OpenAI's AI models like GPT, Codex, and DALLE on Microsoft Azure's cloud platform.

Enterprise Applications: Designed for business use, helping companies integrate AI into various functions such as customer support and content generation.

Scalable Infrastructure: Utilizes Azure's cloud infrastructure for scalable AI model deployment and management.

Security and Compliance: Focuses on high security and adherence to regulatory standards to protect user data.

DeveloperFriendly: Offers tools and APIs for developers to easily implement and customize AI solutions.

Customization Options: Allows users to tailor AI models to better fit their specific requirements.

Global Reach: Available worldwide through Azure's extensive global infrastructure.

Collaboration and Innovation: Supports collaborative AI development and innovation within the Azure ecosystem.

Imp. ChatGPT & OpenAI subscriptions are different

AZURE + Open AI

- collaboration b/w Azure & Open AI
- Secure, Scalable
- It is enterprise version of Open AI

Imp. Not every model is available in every region because of load balancing.

Limit Type	Typical Numbers / Metrics	Notes
Request Rate / Throttling	20–60 requests per minute per deployment (depends on SKU)	Exceeding → 429 errors
Concurrent Requests	5–10 concurrent requests per deployment by default	Exceeding → throttling; can request increase
Token Limits	GPT-4: 8k or 32k tokens, GPT-3.5: 4k tokens	Input + output combined
Deployment Limits	Default ~10 deployments per subscription	Can request more if needed

💡 Interview tip: Mention "Azure OpenAI has request rate limits (20–60/min), concurrent request limits (5–10), and token limits (4k–32k depending on model)"—this shows practical understanding without memorizing every SKU.

Azure Open AI Pricing

GPT-4.1 series

GPT-4.1 series is a highly advanced general-purpose model with extensive world knowledge and an enhanced ability to understand user intent, making it particularly adept at creative tasks and agentic planning. The series features a 1 million token context window and has a knowledge cutoff of June 2024.

Model	Pricing (1M Tokens)	Pricing with Batch API (1M Tokens)
GPT-4.1-2025-04-14 Global	Input: \$2 Cached Input: \$0.50 Output: \$8	Input: \$1 Output: \$4

Model	Pricing
o3-deep-research Global	Input: \$10 Cached Input: \$2.50 Output: \$40

Model	Pricing (1M Tokens)
GPT-5 2025-08-07 Global	Input: \$1.25 Cached Input: \$0.13 Output: \$10

Sign up to Azure Open AI

Microsoft Azure

Home > AI Foundry | Azure OpenAI > Create Azure OpenAI

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 understandings, 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: QBSS1

Resource group: [redacted] [Create new](#)

Instance Details

Region: East US

Name: OpenAI Nad01

Pricing tier: Standard S0

[View full pricing details](#)

Imp. Old = Azure AI Studio
 New = Azure AI Foundry } Azure AI Studio is now AI Foundry

Links

Azure Portal = [portal.azure](https://portal.azure.com)
 Azure AI = [ai.azure](https://ai.azure.com)

Azure AI Foundry → previously AI studio.

Assistants playground
 Speed up development of GPT-powered AI Assistants with prebuilt conversation state management and customization tools.
[Try it now](#) → [Test, Develop playground](#)

Chat playground
 Design a customized AI assistant using ChatGPT. Experiment with GPT-3.5-Turbo and GPT-4 models.
[Try it now](#) → [Simple chat GPT](#)

Bring your own data
 Ground your own data on advanced AI models to create conversational copilot that aid user comprehension, task completion, and decision-making.
[Try it now](#) → [RAG](#)

Images playground
 Generate unique images by writing descriptions in natural language.
[Try it now](#) → [Image Data](#)

Fine-tuning
 Create a custom model by training it with your own data.
[Try it now](#)

AZURE AI MODELS

chat	→ GPT
code	→ Codex
images	→ DALL-E
video	→ SONA
embedding	→ Text-Embedding

↳ basic model
 Then Finetune

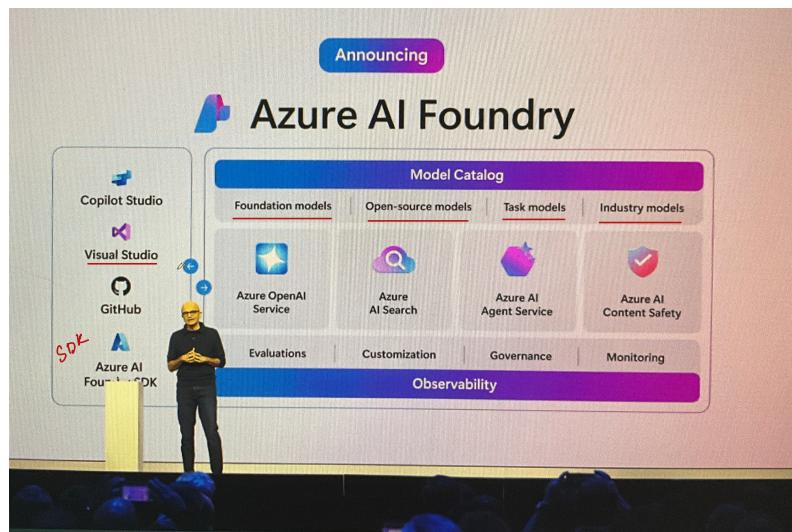
Foundry's Factory
 where metal is
 melted and reshaped

Q: what is Azure AI Foundry?

Azure AI Foundry is Microsoft's all-in-one platform for building, testing, and deploying AI apps and agents. It combines:

- A catalog of AI models (from Microsoft, OpenAI, Meta, etc.) *Imp. Not only meta.*
- Tools for fine-tuning, RAG, and agent creation
- Built-in governance, security, and monitoring
- Support for enterprise use (projects, access control, compliance)

It's designed to help companies go from prototyping to production faster with a unified, scalable AI development environment.



Q: Hubs vs projects?

- Hub: A central workspace for managing **multiple projects**. It's where organizational settings, resources, billing, and policies are set.
- Project: An **isolated workspace** within a hub for building and managing a specific AI solution (e.g., a chatbot, RAG app). Each project has its own data, models, configs, etc.

👉 Think of a **Hub** as the organization or team-level container, and **Projects** as individual AI use cases inside it.

Hub = workspace

Project = individual AI solution

Imp. Azure AI Foundry ; 1. All models
2. Azure Open AI (only openAI)

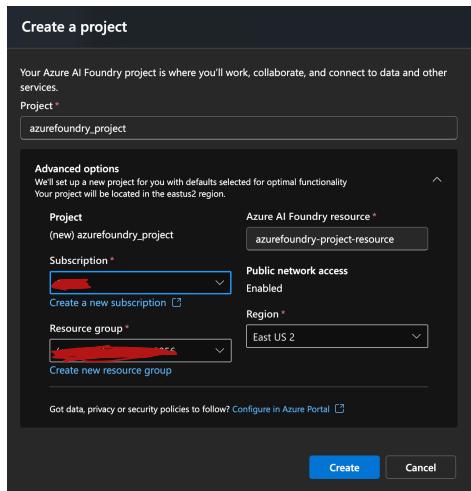
OPTIONS

1. Azure Open AI with project - only Open AI model
2. Azure Open AI without project - other models, SDK

Imp. If you are solely interested in Azure Open AI then there is no need to create a project.

CREATING PROJECT

Go to Azure AI Services



Or how to compose
models?

model catalog



Compose models

DEPLOYING MODELS

Model catalog → Choose model →

Open in
Playground

← Deploying
models

Imp. O₁ are Research Models = more complex /
costly

How to Use

In order to use any of these we need to create deployments



AZURE
AI
SERVICES

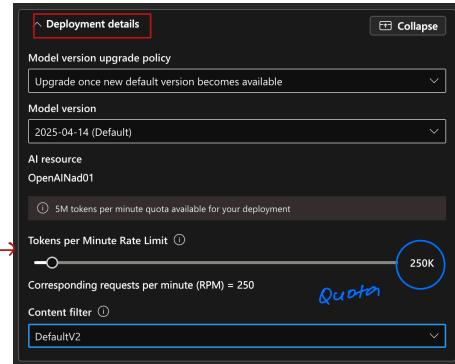
Service	One-liner
Chat	No-code playground to prototype prompts and get responses via Chat Completions before moving to code. learn.microsoft
Assistants	Build configurable copilots that use models plus tools, threads, and files to run tasks. learn.microsoft
Video (Preview)	Playground/APIs to work with multimodal models that understand or generate video content. learn.microsoft
Audio (Preview)	Use speech models (e.g., Whisper family) for transcription, translation, or audio tasks. learn.microsoft
Images	Try image generation models and export code for image workflows. microsoft
Fine-tuning	Train custom versions of base models on domain data and deploy them securely. learn.microsoft
Azure OpenAI Evaluation (Preview)	Evaluate prompts and systems with built-in metrics and datasets to compare quality and safety. learn.microsoft
Stored completions (Preview)	Persist chat/completion histories as datasets to reuse for evals or fine-tuning. learn.microsoft
Batch jobs	Submit large, asynchronous request batches for cost-efficient high-volume processing. learn.microsoft
Monitoring	Track requests, tokens, PTU usage, and fine-tune metrics with Azure Monitor dashboards. learn.microsoft
Deployments	Manage model deployments tied to an endpoint for use in apps and playgrounds. learn.microsoft
Quota	View and manage capacity limits, rate limits, and PTU allocations for the resource. learn.microsoft
Guardrails + Controls	Configure safety filters, content policies, and governance settings for responsible use. learn.microsoft
Risks + alerts (Preview)	Get proactive risk insights and alerts across usage, safety, and operations. learn.microsoft
Data files	Upload and manage files for assistants, fine-tuning, evals, and tool use. learn.microsoft
Assistant	Workspace to configure an individual Assistant's instructions, tools, and files. learn.microsoft
Vector stores (Preview)	Create and attach retrieval indexes for RAG so assistants can ground answers on documents. learn.microsoft

→ Chat Playground

Deploy → Choose model

Imp. If Token limit = 250K

e.g., 2 Deployments = 125K/each



Imp. we don't need to call API here, no code is required

PARAMETERS

Parameter	Meaning (One-liner)
Past Messages	Number of previous conversation turns the model considers for context.
Max Token Completion	Maximum length of the model's response (in tokens).
Temperature	Controls randomness; higher = more creative, lower = more focused.
Top P	↳ 90% probability
Top K	↳ K=2
Presence Penalty	Encourages the model to introduce new topics instead of repeating.
Frequency Penalty	Reduces repetition of words/tokens already used.

Parameter	Explanation	Example
Top-P (Nucleus Sampling)	The model considers the smallest set of tokens whose cumulative probability ≥ P (e.g., P=0.9 → choose from tokens covering 90% of the probability mass).	If possible next words have probs ("cat" 0.5, "dog" 0.3, "fish" 0.1, "car" 0.1), with P=0.9, it will only consider (cat, dog, fish).
Top-K	The model only considers the K most likely tokens and ignores the rest, no matter their probability mass.	With K=2, only ("cat", "dog") are considered even if "fish" had 0.1 probability.

Q: What is use of Azure + OpenAI?

- Enhanced Integration & Enterprise level security.

OpenAI API Calls Vs Azure OpenAI

All models are based on R.L & human Feedback

OpenAI	Azure OpenAI
<p>Python</p> <pre>import os from openai import OpenAI client = OpenAI(api_key=os.getenv("OPENAI_API_KEY"))</pre> <p style="text-align: center;">Open AI</p>	<p>Python</p> <pre>import os from openai import AzureOpenAI client = AzureOpenAI(api_key=os.getenv("AZURE_OPENAI_API_KEY"), api_version="2023-12-01-preview", azure_endpoint=os.getenv("AZURE_OPENAI_ENDPOINT"))</pre> <p style="text-align: center;">Azure OpenAI</p> <p>Store API key in .env file so it is not accessible</p> <p>So 2 additions are; api-version azure Endpoint</p>
<p>OpenAI</p> <p>Python</p> <pre>completion = client.completions.create(model="gpt-3.5-turbo-instruct", prompt=<prompt>) chat_completion = client.chat.completions.create(model="gpt-4", messages=<messages>) embedding = client.embeddings.create(model="text-embedding-ada-002", input=<input>)</pre>	<p>Azure OpenAI</p> <p>Python</p> <pre>completion = client.completions.create(model="gpt-35-turbo-instruct", # This must match the custom deployment name you chose for your model. prompt=<prompt>) chat_completion = client.chat.completions.create(model="gpt-35-turbo", # model = "deployment_name". messages=<messages>) embedding = client.embeddings.create(model="text-embedding-ada-002", # model = "deployment_name". input=<input>)</pre> <p>↳ deployment / not model</p>

Understanding Azure OpenAI

API Calls

```
azure_endpoint = "https://cloudalchemyoai.openai.azure.com/",  
api_key=os.environ.get("OPENAI_API_KEY"),  
api_version="20240215preview"
```

```
completion = client.chat.completions.create(  
    model="air_text_lab", # model = "deployment_name"  
    messages = message_text,  
    temperature=0.7,  
    max_tokens=800,  
    top_p=0.95,  
    frequency_penalty=0,  
    presence_penalty=0,  
    stop=None  
)
```

```
1 #Note: The openai-python library support for Azure OpenAI is in preview.  
2 #Note: This code sample requires OpenAI Python library version 1.0.0 or higher.  
3 import os  
4 from openai import AzureOpenAI  
5  
6 client = AzureOpenAI(  
7     azure_endpoint = "https://cloud-alchemy-oai.openai.azure.com/",  
8     api_key=os.getenv("AZURE_OPENAI_KEY"),  
9     api_version="2024-02-15-preview"  
10 )  
11  
12 message_text = [{"role": "system", "content": "You are an AI assistant that helps people find information."}, {"role": "user", "content": "who is prime minister of India"}, {"role": "assistant", "content": "As of my last update, the Prime Minister of India is Narendra Modi. He has been in office since May 2014. Please note that political positions can change, so it's always a good idea to verify with the latest sources."}]  
13  
14 completion = client.chat.completions.create(  
15     model="air_text_lab", # model = "deployment_name"  
16     messages = message_text,  
17     temperature=0.7,  
18     max_tokens=800,  
19     top_p=0.95,  
20     frequency_penalty=0,  
21     presence_penalty=0,  
22     stop=None  
23 )
```

The screenshot shows the Azure AI services portal with the following details:

- Navigation:** Home > Azure AI services
- Title:** Azure AI services | Azure OpenAI
- Actions:** Create, Manage deleted resources, Manage view, Refresh, Export to CSV, Open query, Assign tags, Delete.
- Filters:** Subscription equals all, Type equals all, Resource group equals all, Location equals all, Add filter.
- Grouping:** No grouping, List view.
- Table Headers:** Name, Kind, Location, Custom Domain Name, Pricing tier, Status, Created date.
- Table Data:**

Name	Kind	Location	Custom Domain Name	Pricing tier	Status	Created date
cloud-alchemy-oai	OpenAI	East US	cloud-alchemy-oai	S0	Succeeded	2024-05-01T13:25:43.685Z
- Bottom Navigation:** Overview, All Azure AI services, Azure AI services, Azure OpenAI, AI Search.

AZURE Endpoint URL & API Keys?

These are very important for making API calls

The screenshot shows the Azure portal interface for the 'OpenAINad01' resource group. The left sidebar has 'Keys and Endpoint' selected under 'Resource Management'. The main pane displays two API keys ('KEY 1' and 'KEY 2'), their locations ('eastus'), and the endpoint URL ('https://openainad01.openai.azure.com/'). A note at the top of the pane advises users to store keys securely.

Step 1: AZURE Open AI API Calling

Unlike OpenAI call,

we need

- Azure AI API key
- API version
- Client ID

The screenshot shows the 'Welcome to Azure OpenAI' configuration page. It displays the 'Resource configuration' section with fields for 'Name' (OpenAINad01), 'Subscription' (selected), 'Subscription ID' (redacted), 'View access control (IAM)', 'API key 1' (redacted), 'Resource group' (naditest001), 'Pricing tier' (Standard S0), 'Azure OpenAI endpoint' (https://openainad01.openai.azure.com/), and 'Location' (eastus). The 'API key 1' field and the endpoint URL are circled in yellow.

Q: How to get code;

1. Set up Key, endpoint, in .env file

→ Deploy the chat model (use it in client)

Imp:

Code = Chatplay Ground → View Code

```

import os
from dotenv import load_dotenv
load_dotenv()
from openai import AzureOpenAI

AZURE_API_CLIENT = os.getenv("AZURE_API_CLIENT")
AZURE_API_KEY = os.getenv("AZURE_API_KEY")

client= AzureOpenAI(
    api_key=AZURE_API_KEY,
    api_version= "2025-01-01-preview",
    azure_endpoint=AZURE_API_CLIENT)
  
```

```

chat_prompt = [
    {"role": "user",
     "content": "Who has won the most FIFA World Cups in soccer?"}
]

# Generate the response

response = client.chat.completions.create(
    model="my-first-gpt",
    messages=chat_prompt
)

print(response.choices[0].message.content)
  
```



History behind Azure OpenAI

- **2019 Partnership:** Microsoft and OpenAI formed a partnership, emphasizing the integration of OpenAI's technologies with Microsoft Azure.
- **\$1 Billion Investment:** Microsoft announced an investment of \$1 billion into OpenAI as part of their partnership agreement.
- **2021 Azure OpenAI Service Launch:** Launched to enterprise customers, allowing them access to OpenAI's AI models like GPT3, through Microsoft's Azure platform.
- **Expanding Capabilities:** The service has since expanded to include more sophisticated AI models, such as Codex and DALLE, enhancing its offerings for complex enterprise applications.
- **Enterprise Focus:** Designed specifically for enterprise use, Azure OpenAI aims to support businesses in incorporating AI into their operations securely and at scale.