# Generative AI-1

Azure for Students
azuregenai

← All workspaces

🏠 Home
🗂 Model catalog

**Authoring**

📓 Notebooks
⚡ Automated ML
🔧 Designer
>_ Prompt flow
🔬 Tracing PREVIEW

**Assets**

📊 Data
🧪 Jobs
🗃 Components
🔀 Pipelines

# azuregenai ✎

+ New    ⚙ Customize view

## Generative AI with Prompt flow ⋯

View prompt flow →

### Multi-Round Q&A on Your Data

Create a chatbot that uses LLM and data from your own indexed files to ground multi-round question and answering capabilities in enterprise chat scenarios.

Start    Clone

### Q&A on Your Data

Use LLM and data from your own indexed files to ground multi-round question and answering capabilities.

Start    Clone

### Web Classification

Use LLM to classify URLs into multiple categories.

Start    Clone

### Chat with Wikipedia

Create a chatbot that leverages Wi to ground the responses.

Start    Clone

## Generative AI models ⋯

View all →

🔷 gpt-4o          ✅

🔷 gpt-4          ✅

Ⓜ mistralai-Mixtral-8x7B-v01    ✅

Ⓜ mistralai-Mixtral-8x7B-    ✅

---

Azure for Students
azuregenai

← All workspaces

🏠 Home
🗂 Model catalog

**Authoring**

📓 Notebooks
⚡ Automated ML
🔧 Designer
>_ Prompt flow
🔬 Tracing PREVIEW

**Assets**

📊 Data
🧪 Jobs
🗃 Components
🔀 Pipelines

Default Directory > azuregenai > Model catalog

## Find the right model to build your custom AI solution

⊞ Show filters

🟦 Phi-3-medium-128k-instruct    ✅
Chat completion

🟦 Phi-3-vision-128k-instruct    ✅
Chat completion

∞ Meta-Llama-3-8B-Instruct    ✅
Chat completion

∞ Meta-Llama-3-8B    ✅
Text generation

∞ Meta-Llama-3-70B-Instruct    ✅
Chat completion

∞ Meta-Llama-3-70B    ✅
Text generation

∞ Llama-2-7b-chat    ✅
Chat completion

∞ Llama-2-7b    ✅
Text generation

∞ Llama-2-70b-chat    ✅
Chat completion

∞ Llama-2-70b    ✅
Text generation

∞ Llama-2-13b-chat    ✅
Chat completion

∞ Llama-2-13b    ✅
Text generation

∞ CodeLlama-7b-hf    ✅
Text generation

∞ CodeLlama-7b-Python-hf    ✅
Text generation

∞ CodeLlama-34b-hf    ✅
Text generation

∞ CodeLlama-34b-Python-hf    ✅

∞ CodeLlama-34b-Instruct-hf    ✅

∞ CodeLlama-13b-hf    ✅

Default Directory > azuregenai > Model catalog > azureml-meta > Meta-Llama-3-70B-Instruct

# Meta-Llama-3-70B-Instruct ✓ PREVIEW

Overview    Versions    Artifacts    Security

🔘 Task: Chat completion    💉 Fine-tuning task: chat-completion    🌐 Languages: EN    🏢 License: custom

↻ Refresh    📈 Evaluate    ⚗ Fine-tune    ▷ Deploy    📜 View license

## Description

### Model Details

Meta developed and released the Meta Llama 3 family of large language models (LLMs), a collection of pretrained and instruction tuned generative text models in 8 and 70B sizes. The Llama 3 instruction tuned models are optimized for dialogue use cases and outperform many of the available open source chat models on common industry benchmarks. Further, in developing these models, we took great care to optimize helpfulness and safety.

**Model developers** Meta

**Variations** Llama 3 comes in two sizes — 8B and 70B parameters — in pre-trained

### Serverless APIs PREVIEW

Provision an inference API within seconds through Models as a Service. Explore the model in the playground. Pay only for tokens consumed.

**Pricing**

| paygo-inference-output-tokens: $0.01134 per 1000 tokens | paygo-inference-input-tokens: $0.00378 per 1000 tokens |

---

## Deployment options

| Serverless API with Azure AI Content Safety PREVIEW 🗓 | Managed Compute without Azure AI Content Safety ⚠ |
|---|---|
| This option offers managed API service that does not require you to host or manage infrastructure. You can choose to include standard Azure AI Content Safety filters with this option. | This option offers user-managed hosting and model inferencing on Azure infrastructure. It doesn't use Azure AI Content Safety filters and you may be at higher risk of exposing users to harmful content. |

industry benchmarks. Further, in developing these models, we took great care to optimize helpfulness and safety.

**Model developers** Meta

**Variations** Llama 3 comes in two sizes — 8B and 70B parameters — in pre-trained

| paygo-inference-output-tokens: $0.01134 per 1000 tokens | paygo-inference-input-tokens: $0.00378 per 1000 tokens |

## Azure AI | Machine Learning Studio

Azure for Students
azuregenai

All workspaces

Home

Model catalog

**Authoring**

Notebooks

Automated ML

Designer

Prompt flow

Tracing PREVIEW

**Assets**

Data

Jobs

Components

Pipelines

Default Directo

Meta-Lla

Overview

Task: Chat c

Refresh

Description

Model Det

Meta develope
(LLMs), a colle
and 70B sizes.
cases and outp
industry bench
optimize helpfu

Model develop

Variations Llai

**Serverless API deployment for Meta-Llama-3-70B-Instruct**

PREVIEW

×

Overview    Pricing and terms

Meta

Meta Llama-3-70B Instruct is offered by Meta AI through the Azure Marketplace. View the pricing and terms tab to learn about pricing and terms of use.
Learn more about Models as a Service. ⧉

**Terms of use**

By clicking "Subscribe and Deploy", I (a) agree to the legal terms and privacy statements associated with each Marketplace offering above, (b) authorize Microsoft to charge or bill my current payment method for the fees associated with my use of the offerings, including applicable taxes, with the same billing frequency as my Azure subscription, until I discontinue use of the offerings, (c) agree that Microsoft may share my contact information and transaction details (including usage volume associated with the offering) with the sellers of the offerings so that they can contact me regarding this product. Microsoft does not provide rights for third

Azure Marketplace Terms ⧉

seconds through Models as a Service. Explore only for tokens consumed.

paygo-inference-input-tokens:
$0.00378 per 1000 tokens

**Subscribe and Deploy**    Cancel

---

ap4ashutosh / Azure-GenAI-Playground

Type / to search

<> Code    ⊙ Issues    ⁑ Pull requests    ⊙ Actions    ⊞ Projects    ⊙ Security    ⌁ Insights

**Azure-GenAI-Playground**  Public

👁 Watch  0 ⌄    ⑂ Fork  19 ⌄    ☆ Star  0 ⌄

main    ⑂ 1 Branch    ⬡ 0 Tags

Go to file    t    +    <> Code ⌄

**About**

*No description, website, or topics provided.*

ap4ashutosh  updation                          35f1786 · 10 hours ago    ⊙ 2 Commits

📁 project1          updation                              10 hours ago
📄 .env             updation                              10 hours ago
📄 LICENSE          commiting 1st project                 yesterday
📄 README.md        commiting 1st project                 yesterday
📄 test.py          updation                              10 hours ago

📖 README    ⚖ GPL-2.0 license

📖 Readme

⚖ GPL-2.0 license

⌁ Activity

☆ 0 stars

👁 0 watching

⑂ 19 forks

Report repository

**Releases**

No releases published

```
SUBHAM@LAPTOP-D67N4JI9 MINGW64 ~
$ cd /e

SUBHAM@LAPTOP-D67N4JI9 MINGW64 /e
$ git clone https://github.com/ap4ashutosh/Azure-GenAI-Playground
```

File   Edit   Selection   View   Go   ⋯                                    🔍 project1                                    ◻ ▭ ◫ ▱    ─ ▢ ✕

EXPLORER                              ⋯         ● chat_be.py  ✕                                                                          ▷ ▯ ⋯

∨ PROJECT1                                       chat > ● chat_be.py

∨ 🗀 chat                                    1    # This is the backend file for the chat application using models from Azure ML model catalog
   ⋗ 🗀 __pycache__                         2    # Written by Ashutosh
     ⚙ .env                                3    # Date: 15th june 2024
   ● chat_be.py                            4
   ● chat_fe.py                            5    import os
   📄 req.txt                              6    from dotenv import load_dotenv, find_dotenv
   📄 text_completion.py                   7    from azure.ai.ml import MLClient
                                           8    from azure.identity import DefaultAzureCredential
                                           9
                                          10    from langchain_community.chat_models.azureml_endpoint import AzureMLChatOnlineEndpoint, CustomOpenAIChatContentFormatter
                                          11    from langchain_core.messages import HumanMessage
                                          12    from langchain.memory import ConversationBufferMemory
                                          13    from langchain.chains import ConversationChain
                                          14
                                          15
                                          16    # Load environment variables from .env file
                                          17    load_dotenv(find_dotenv())
                                          18
                                          19    # Fetch values from environment variables
                                          20    subscription_id = os.getenv("subscription_id")
                                          21    resource_group = os.getenv("resource_group")
                                          22    workspace = os.getenv("workspace")
                                          23    url = os.getenv("endpoint_url")
                                          24    api = os.getenv("endpoint_api_key")
                                          25

PROBLEMS   OUTPUT   DEBUG CONSOLE   TERMINAL   PORTS                                        ⊘ powershell  + ⌄  ▯  🗑  ⋯  ⌃  ✕

⊘ PS E:\Azure-GenAI-Playground\project1>

∨ OUTLINE
∨ TIMELINE

⊗ 0 △ 0   ⓦ 0                                                                  Ln 1, Col 1   Spaces: 4   UTF-8   CRLF   Python   ☐