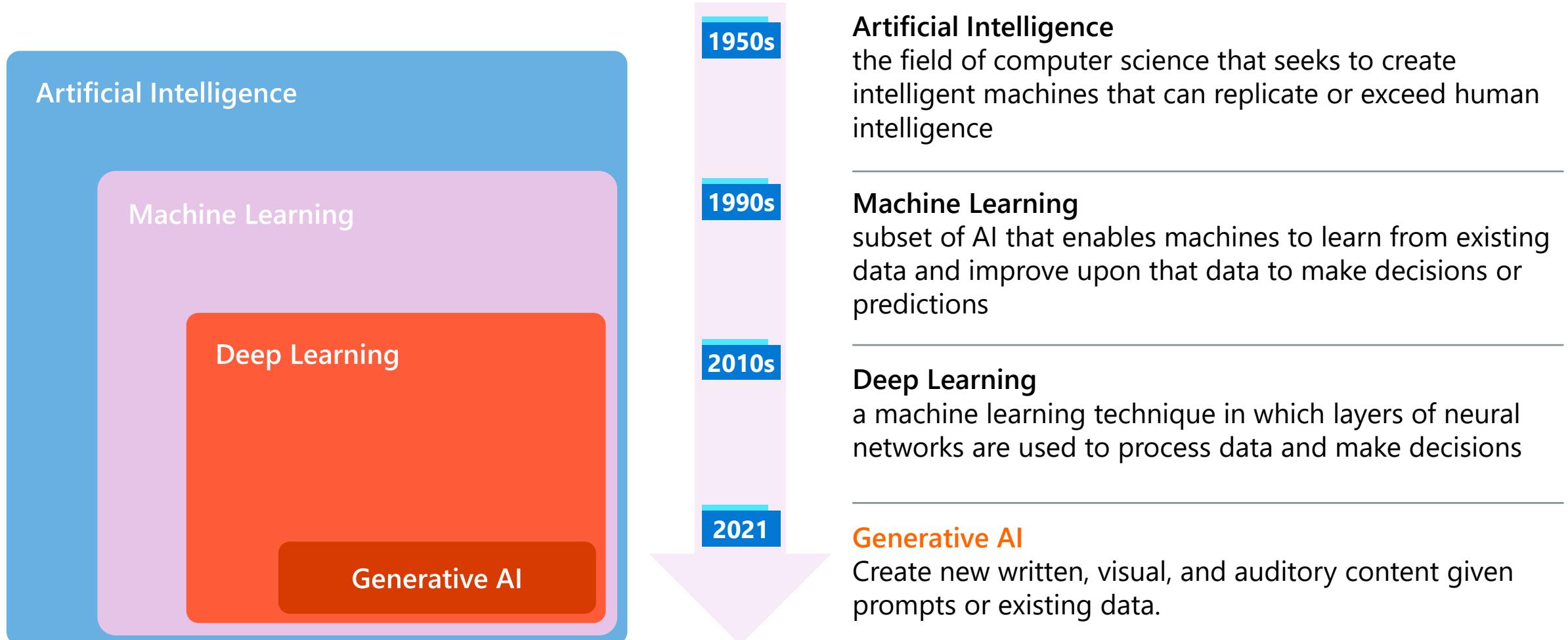


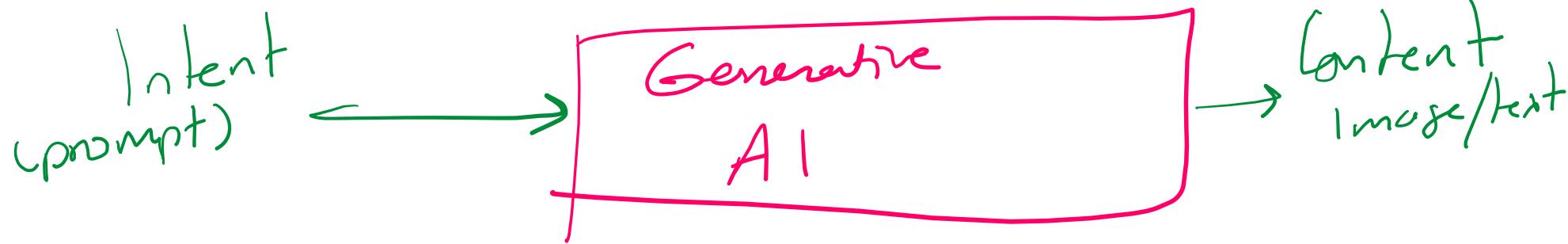
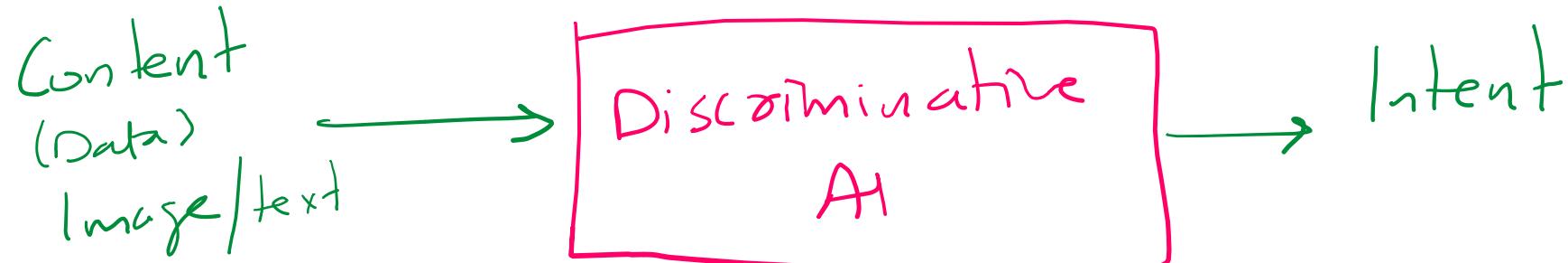


# Generative AI & Prompt Engineering with Azure OpenAI



# What is generative AI?





# Tasks in Generative AI: Text Generation

- **Text Generation:** Creating human-like text based on input prompts.
- Latest models/platforms: OpenAI's GPT-4, Google's PaLM & Gemini, Meta's LLaMA.

The screenshot shows a dark-themed interface of the ChatGPT web application. At the top right, it says "Model: GPT-4". A user message is shown with a green square icon containing a white "B" and the text: "How can I use storytelling techniques to make my videos more compelling?". Below it, the AI response is shown with a blue square icon containing a white "G" and the text: "Using storytelling techniques in your videos can make them more engaging, memorable, and compelling for your audience. Here are some key storytelling elements and strategies to consider when creating your videos: 1. [redacted]". At the bottom left, there is a text input field with the placeholder "Send a message." and a "Stop generating" button with a checkbox. At the very bottom, a small note reads "ChatGPT may produce inaccurate information about people, places, or facts. ChatGPT May 3 Version".

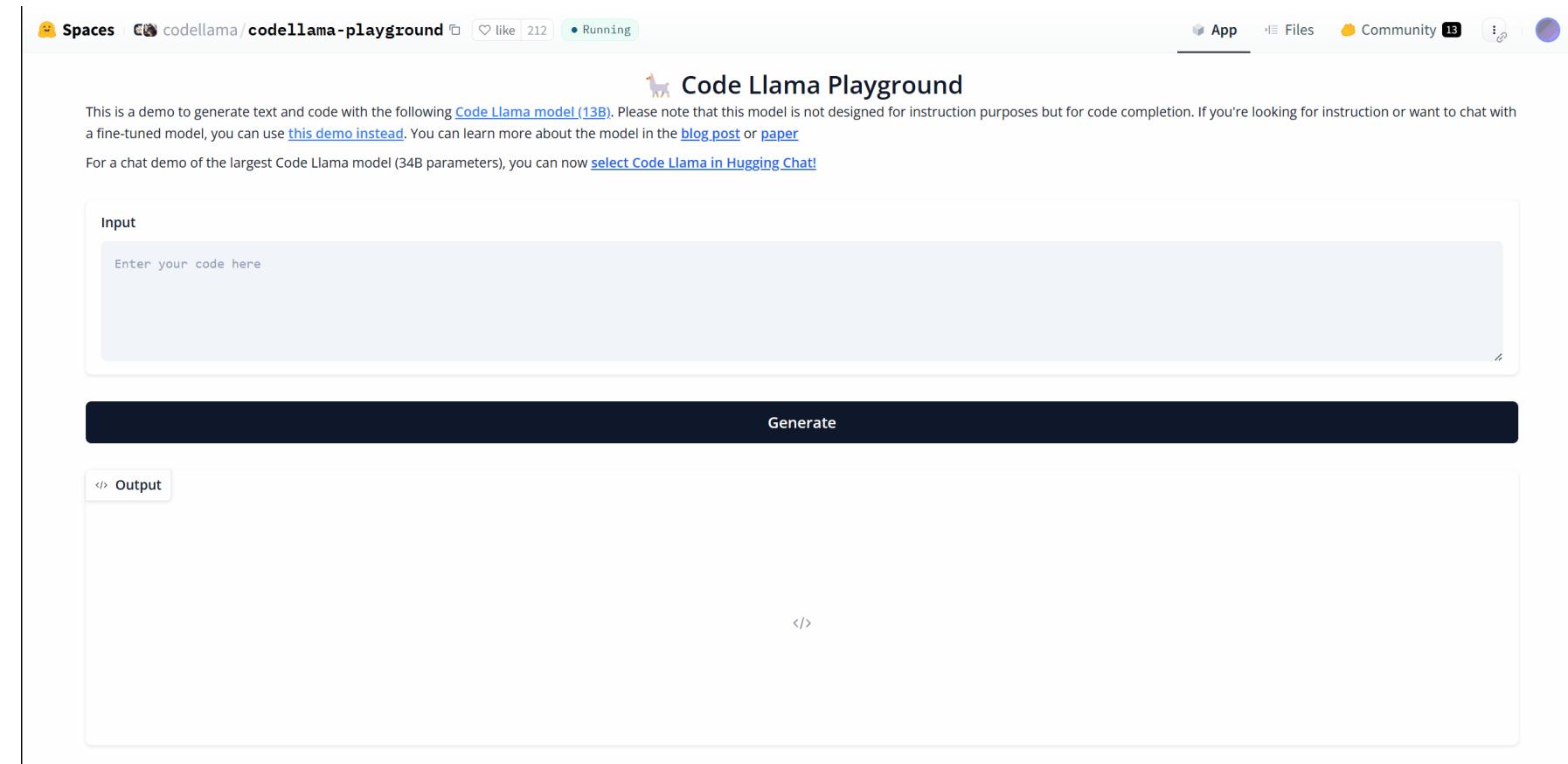
# Tasks in Generative AI: Text Generation

Models like GPT-4-vision, llava-v1 can also be used to generate text from images, taking images as prompt.

The screenshot shows the LLaVA interface on a web browser. The top navigation bar includes 'Spaces' (badayvedat/LLaVA), 'like 325', 'Running on T4', 'App', 'Files', 'Community 8', and user profile icons. The main title is 'LLaVA: Large Language and Vision Assistant' with links to [Project Page], [Paper], [Code], and [Model]. A note states 'ONLY WORKS WITH GPU!'. It explains that the model can be loaded with 4-bit or 8-bit quantization, and recommends 4-bit for A10G for best efficiency. Recommended configurations are listed for Hardware (T4-Small, A10G-Small, A100-Large) and Bits (4 default, 4, 16). A dropdown menu shows 'llava-v1.5-13b-4bit'. Below it is an 'Image' input field with a placeholder 'Drop Image Here - OR - Click to Upload'. An 'Examples' section shows a thumbnail of a person on a bicycle next to the question 'What is unusual about this image?'. On the right side, there is a 'LLaVA Chatbot' section and a 'Record' button at the bottom right.

# Tasks in Generative AI: Code Generation

- **Code Generation:** Automatically writing or suggesting code snippets based on a user's instructions or existing code context.
- Latest models/platforms: OpenAI's Codex (which powers GitHub Copilot), Google's AlphaCode, Amazon's CodeWhisperer.



# Tasks in Generative AI: Image Generation

- **Image Generation:**  
Generating photorealistic or artistic images from textual descriptions.
- **Latest models/platforms:**  
OpenAI's DALL-E 2, Stability AI's Stable Diffusion, Google's Imagen.

The screenshot shows the homepage of the Stable Diffusion 2.1 Demo. At the top, there is a navigation bar with links for 'Spaces', 'stabilityai/stable-diffusion' (with a fork icon), 'like 10.2k', and 'Running on CPU UPGRADE'. On the right side of the nav bar are links for 'App', 'Files', 'Community 19624', and a profile icon. Below the nav bar, the title 'Stable Diffusion 2.1 Demo' is displayed in bold. A sub-header below it reads: 'Stable Diffusion 2.1 is the latest text-to-image model from StabilityAI. [Access Stable Diffusion 1 Space here](#). For faster generation and API access you can try [DreamStudio Beta](#)'. The main interface features two input fields: 'Enter your prompt' and 'Enter a negative prompt', followed by a large central area for image generation. To the right of these fields is a 'Generate image' button. At the bottom of the page, there is a 'Share to community' button and a link for 'Advanced settings'.

# Tasks in Generative AI: Audio Generation

- **Audio Generation:** Producing synthetic speech, music, or other sound effects.
- Latest models/platforms: Google's AudioLM, OpenAI's Jukebox, Descript's Overdub.

The screenshot shows the MusicGen demo interface within a Streamlit application. The top navigation bar includes links for Spaces, facebook, MusicGen, like 3.92k, and Running on A10G. The main title "MusicGen" is displayed, followed by a brief description: "This is the demo for [MusicGen](#), a simple and controllable model for music generation presented at: ["Simple and Controllable Music Generation"](#)". A "Duplicate Space" button is available for longer sequences, more control, and no queue. The interface features three main input fields: "Describe your music" (with a text area), "Condition on a melody (optional) File or Mic" (with "file" and "mic" radio buttons, a "File" button, and a "Drop Audio Here - OR - Click to Upload" area), and "Generated Music" (with a "Generated Music (wav)" file preview). A large "Generate" button is located below the inputs. At the bottom, there is a section titled "Examples" with three rows of generated music descriptions and corresponding file links:

Describe your music	File
An 80s driving pop song with heavy drums and synth pads in the background	bach.mp3
A cheerful country song with acoustic guitars	bolero_ravel.mp3
90s rock song with electric guitar and heavy drums	

## Tasks in Generative AI: Video Generation

- Video Generation:  
Creating video clips from scratch or modifying existing videos with generative methods.
- Latest models/platforms:  
OpenAI's Sora, Meta's Make-A-Video, Google's DreamFusion, Synthesia for AI video production.



Prompt: A young man at his 20s is sitting on a piece of cloud in the sky, reading a book.

# Large Language Models

Dutz

I love pizza	+ve	$[2 \cdot 1, -3 \cdot 1, 4 \cdot 6, 4 \cdot 2]$	+ve
I hate pasta	-ve	$[3 \cdot 1, 4 \cdot 2, 1 \cdot 9, 6 \cdot 1, \cancel{2 \cdot 2}]$	-ve
I love noodles	+ve	$[4 \cdot 3, 4 \cdot 6, 3 \cdot 1, 2 \cdot 5]$	+ve

A hand-drawn diagram consisting of a blue arrow pointing upwards and to the right, forming a curve that dips down and then rises. The word "Vectorization" is written in blue below the arrow.

*Data*

11

# ML Algorithm

Model

→ trained model

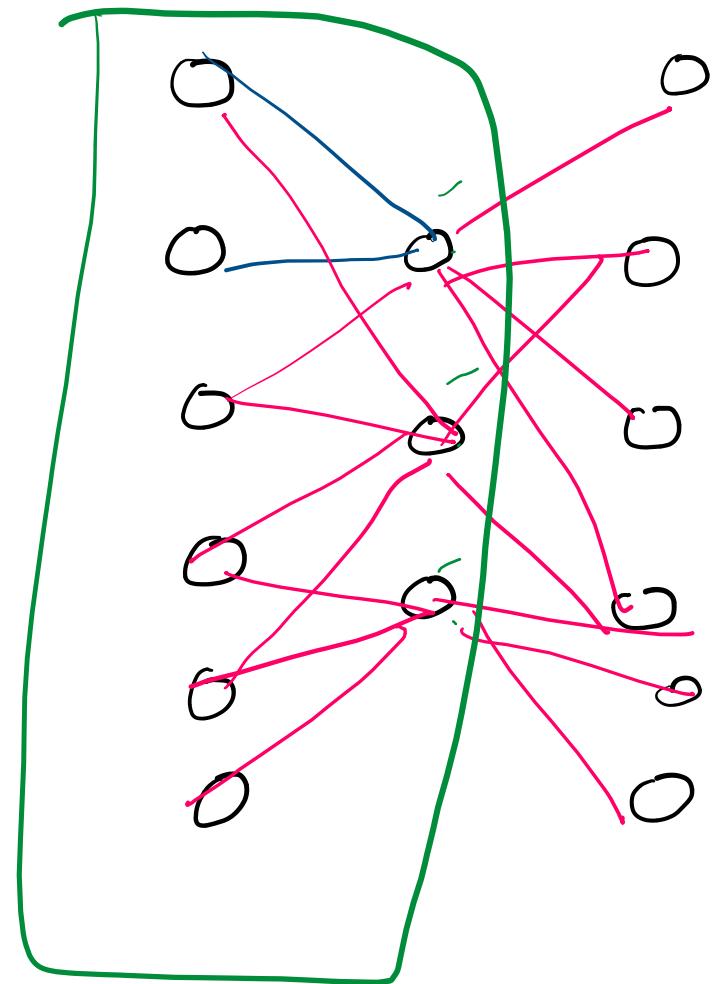
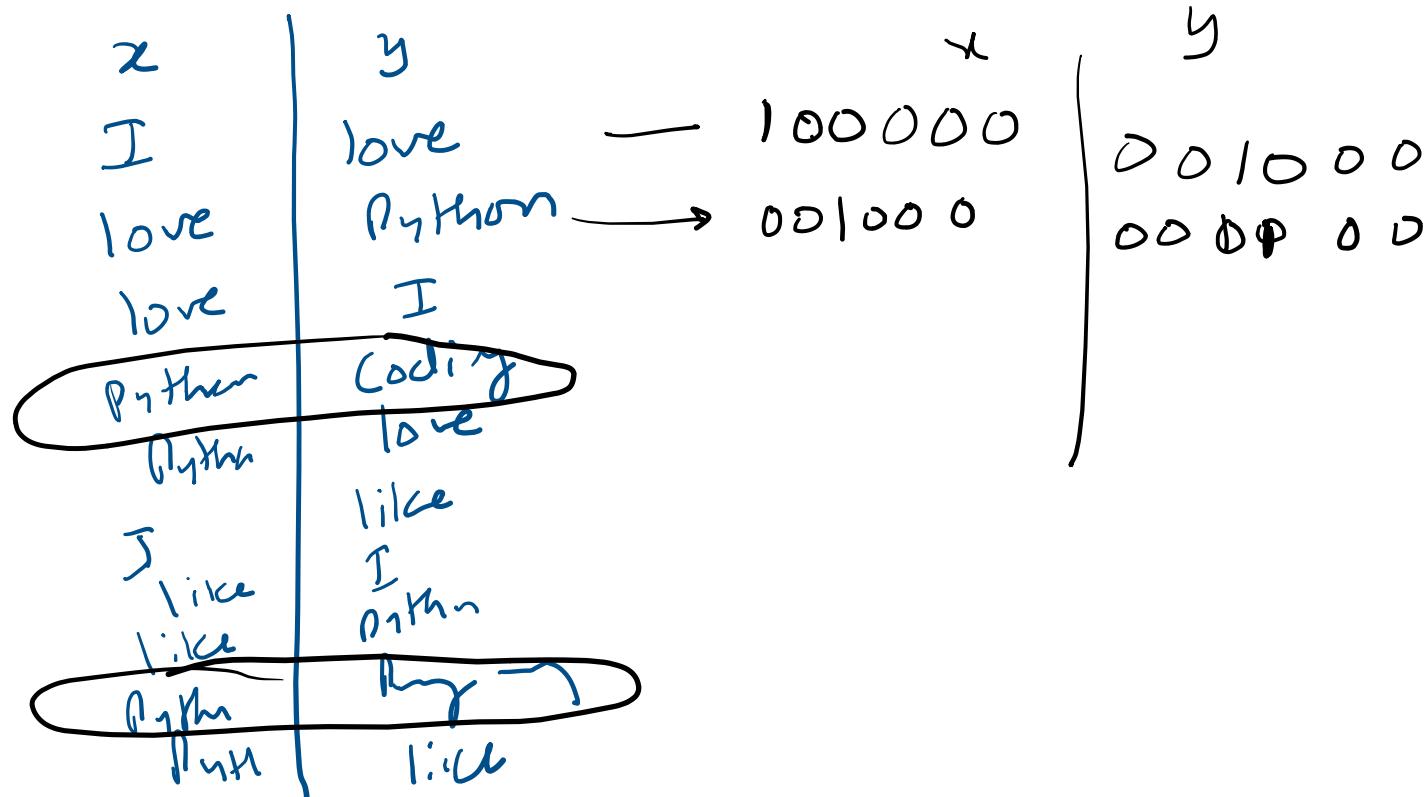
# Large Language Models

(Woodwec)

③ → vector size

- ① I love python coding
  - ② I like python programming.

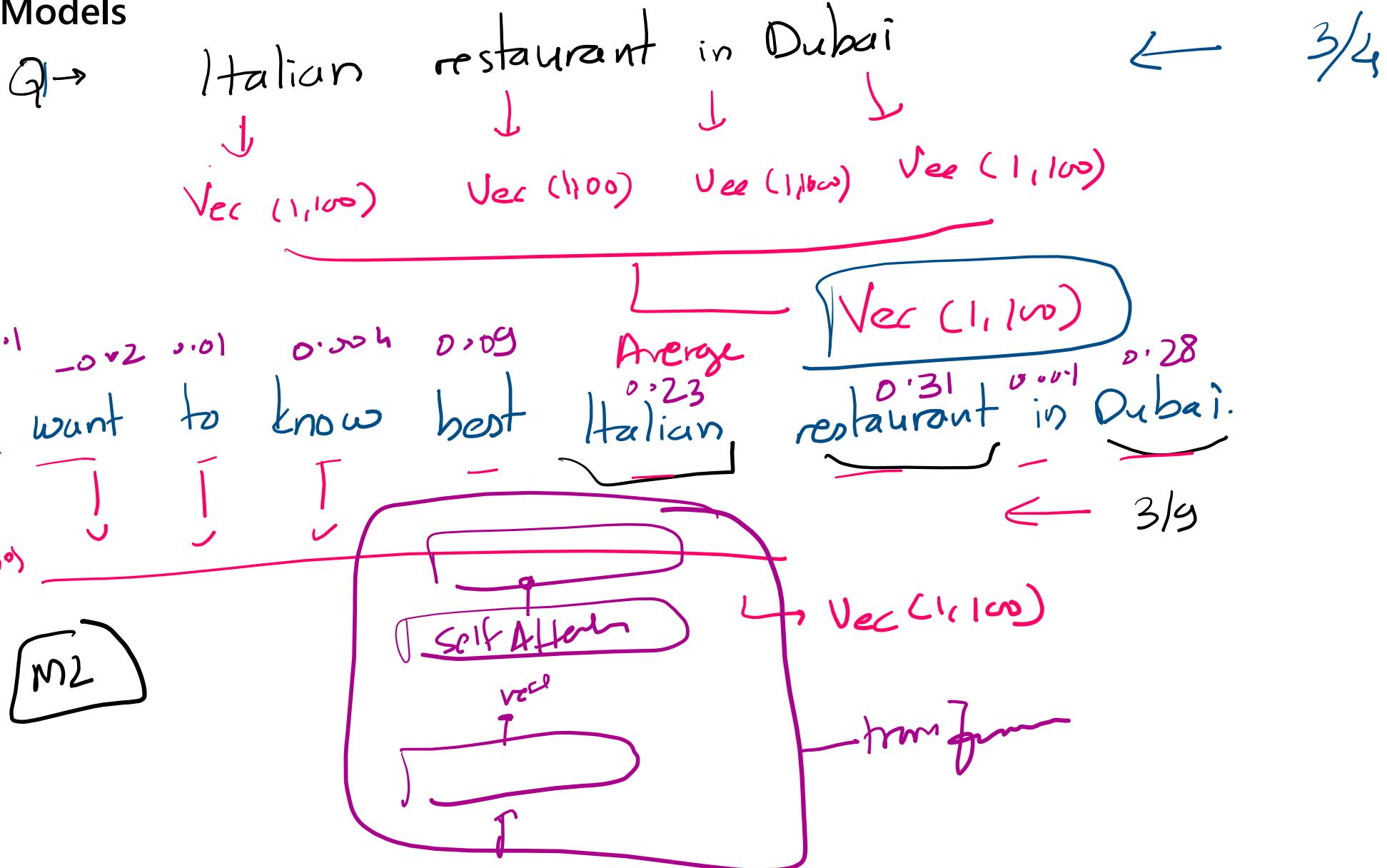
Unique → [ I, like, love, python, Programming . cody ]



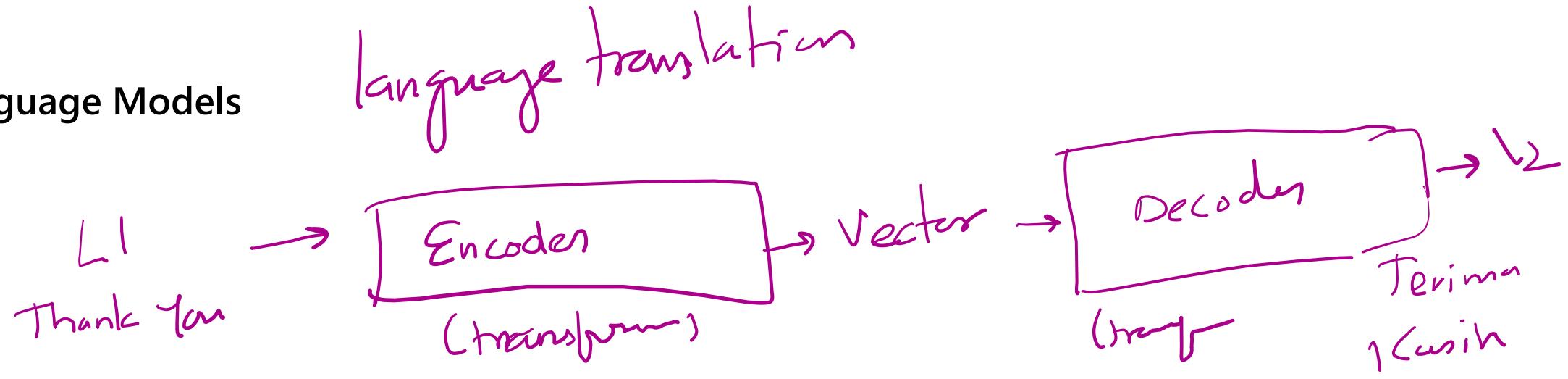
## Vector model

2012 - 2016

## Large Language Models



## Large Language Models



# Large Language Models (LLMs)



Type of AI that can process and produce natural language text.



Learns from a massive amount of data gathered from different sources.



A foundation model refers to a specific instance or version of an LLM.  
Example - GPT-3, GPT-4.

# Azure OpenAI Service

# Azure OpenAI Service



Deployed in your Azure subscription, secured by you, and tied to your datasets and applications



Large, pretrained AI models to unlock new scenarios



AI models, some custom-tunable with your data and hyperparameters



Built-in responsible AI to detect and mitigate harmful use



Enterprise-grade security with role-based access control (RBAC) and private networks

# Top use cases

Business Problem	Productivity is lagging	Need for process Automation	Degraded Customer Experience	Creating Content is Time Consuming
Business Needs	Increase Productivity	Automate Processes	Improve Customer Experience	Build Creative Content
Solutions	<ul style="list-style-type: none"><li>Conversational Search/Knowledge Insights</li><li>Code Generation and Documentation</li><li>Trend Forecasting</li><li>Report Summarization &amp; Generation</li></ul>	<ul style="list-style-type: none"><li>Document Processing</li><li>Workflow Management</li><li>Fraud Detection</li><li>Supply Chain Optimization</li></ul>	<ul style="list-style-type: none"><li>Intelligent Contact Center</li><li>Agent/Employee Assistance</li><li>Virtual Assistance</li><li>Call Analytics</li><li>Call Summarization</li></ul>	<ul style="list-style-type: none"><li>Marketing/Sales Content Generation</li><li>Personalized Content Generation</li><li>Product Design &amp; Development</li><li>Digital Art</li></ul>
What can Generative AI Do?	<p>Generate New Revenue Streams Deliver Differentiated Customer Experiences Modernize Internal Processes</p>			

# Top capabilities and scenarios

1.0



Content generation



Summarization



Code generation



Semantic search

## Examples of advanced use cases

2.0

### Telecommunications

Media Workflows, Cross Content Linking, Content Creation for Media, Speech Analytics, Analytics for B2C Contact Center, Cognitive Contact Center, Skilling Automation, Realtime Speech Transcriptions & Summarization

### Manufacturing & Industrials

ChatGPT Enabled Technical Support, Customer Sentiment Analysis, Customer Service Knowledge Mining, Digital Proposal Assistant, Customer Journey Analytics, Consumer Insights Advanced Analytics, Records Summarization, Anomaly Detection, Virtual Agents with Copilot

### Automotive, Mobility & Transportation

Marketing Content Generation, Contextual Contact Center, Customer Feedback Loop, Smart Incident Manager, Customer Comms, Text Summarization & Analytics

## Customers Bringing it to Life



# Components and key concepts of Azure OpenAI

# Azure Open AI - Components



Pre-trained  
generative AI  
models



Customization  
capabilities



Built-in tools to  
detect and mitigate  
harmful use cases



Enterprise-grade  
security with RBAC

# Key concepts



Prompts & Completions



Tokens



Resources



Deployments



Prompt engineering



Models

# Provision an Azure OpenAI resource in Azure

## Deploy a model in Azure OpenAI Studio to use it

1. Apply for access to the Azure OpenAI service:  
<https://aka.ms/oaiapply>
2. Create an **Azure OpenAI** resource in the Azure portal

Alternatively, use the Azure CLI

```
az cognitiveservices account create \
-n MyOpenAIResource \
-g MyResourceGroup \
-l eastus \
--kind OpenAI \
--sku s0 \
--subscription subscriptionID
```

Home > Azure AI services | Azure OpenAI >

## Create Azure OpenAI ...

1 Basics   2 Network   3 Tags   4 Review + submit

Enable new business solutions with OpenAI's language generation capabilities powered by GPT-3 models. These models have been pretrained with trillions of words and can easily adapt to your scenario with a few short examples provided at inference. Apply them to numerous scenarios, from summarization to content and code generation.

[Learn more](#)

### Project Details

Subscription \* ⓘ

Resource group \* ⓘ

[Create new](#)

### Instance Details

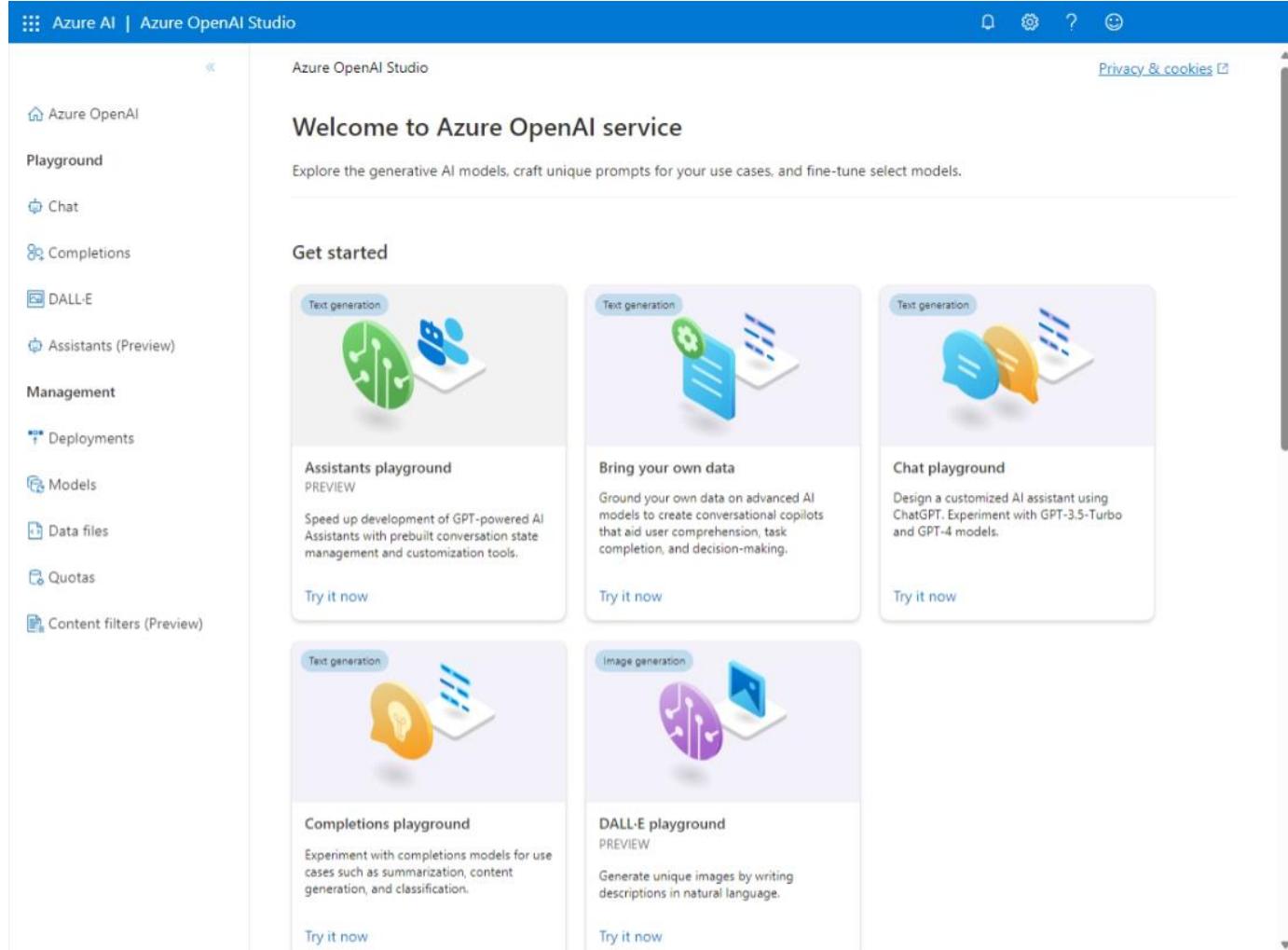
Region ⓘ

Name \* ⓘ

Pricing tier \* ⓘ

# Azure OpenAI Studio

- Web portal for working with Azure OpenAI models:  
<https://oai.azure.com/>
- View and deploy base models
- Connect your own data source
- Manage fine tuning and data files for custom models
- Test models in visual playgrounds:
  - **Chat** (GPT-3.5-Turbo and later models)
  - **Completions** (GPT-3 and earlier models)
  - **DALL-E** (Image generations)
  - **Assistants** (Custom and Copilot-like experiences)



The screenshot shows the Azure OpenAI Studio interface. The left sidebar lists various services: Azure OpenAI, Playground, Chat, Completions, DALL-E, Assistants (Preview), Management, Deployments, Models, Data files, Quotas, and Content filters (Preview). The main area is titled "Welcome to Azure OpenAI service" with the sub-section "Get started". It features six cards arranged in a 2x3 grid:

- Assistants playground** (PREVIEW): Speed up development of GPT-powered AI Assistants with prebuilt conversation state management and customization tools. [Try it now](#)
- Bring your own data**: Ground your own data on advanced AI models to create conversational copilots that aid user comprehension, task completion, and decision-making. [Try it now](#)
- Chat playground**: Design a customized AI assistant using ChatGPT. Experiment with GPT-3.5-Turbo and GPT-4 models. [Try it now](#)
- Completions playground**: Experiment with completions models for use cases such as summarization, content generation, and classification. [Try it now](#)
- DALL-E playground** (PREVIEW): Generate unique images by writing descriptions in natural language. [Try it now](#)

At the top right of the main area, there are links for "Privacy & cookies" and a "Logout" button.

# Types of generative AI model

Model Family	Description
GPT-4	Newest, most capable chat-based models for language and code generation
GPT-3.5	Natural language and code-generation models
Embeddings	Models that use embeddings for specific tasks (similarity, text search, and code search)
DALL-E	Image-generation model

The screenshot shows the Azure AI Studio interface with the 'Models' section selected. The left sidebar includes links for Azure OpenAI, Playground, Chat, Completions, DALL-E (Preview), Management, Deployments, and Models (which is highlighted). The main content area displays a table titled 'Base models' with columns for Model name, Model version, Created at, Status, and Deployable. The table lists several models:

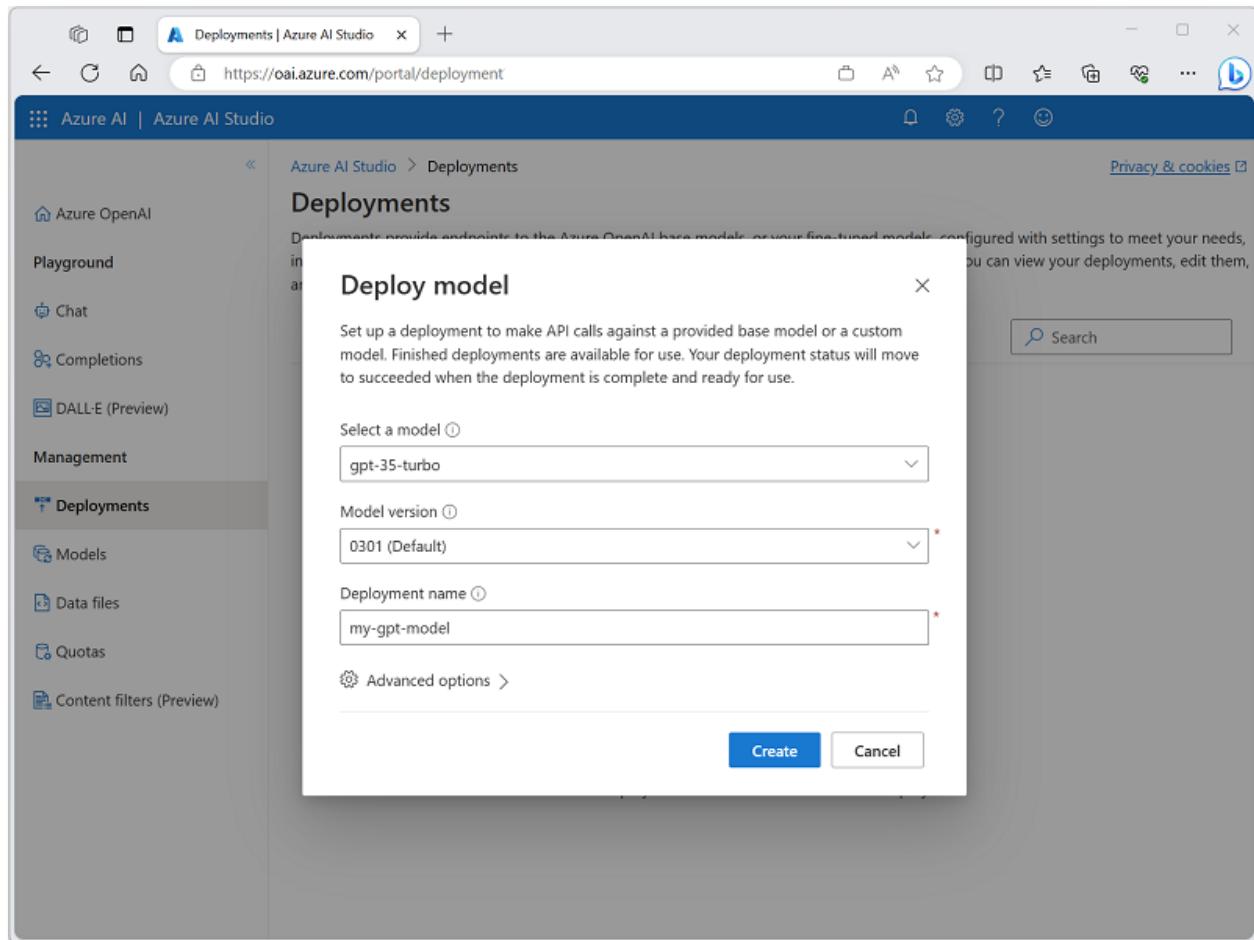
Model name	Model version	Created at	Status	Deployable
gpt-35-turbo	0613	6/18/2023 5:00 PM	Succeeded	Yes
gpt-35-turbo	0301	3/8/2023 4:00 PM	Succeeded	Yes
gpt-35-turbo-16k	0613	6/18/2023 5:00 PM	Succeeded	Yes
text-embedding-ada-002	2	4/2/2023 5:00 PM	Succeeded	Yes
text-embedding-ada-002	1	2/1/2023 4:00 PM	Succeeded	Yes

# Deploying generative AI models

## Deploy a model in Azure OpenAI Studio to use it

- You can deploy one or more instances of each available model
- The number of deployments depends on your quota, which you can see in the portal
- Alternatively, use the Azure CLI

```
az cognitiveservices account deployment create \
  -g myResourceGroupName \
  -n MyOpenAIREsource \
  --deployment-name my-gpt-model \
  --model-name gpt-35-turbo \
  --model-version "0301" \
  --model-format OpenAI \
  --scale-settings-scale-type "Standard"
```





## Using prompts to get completions from models

Task	Prompt	Completion
Classifying content	Tweet: I enjoyed the training course. Sentiment:	Positive
Generating new content	Write a poem about databases	Databases, oh databases, You keep our information safe, From the small to the large, You store our data in a place.
Transformation/Translation	English: Hello French:	Bonjour
Summarization	Scotland is [long description of Scotland...]  Summarize the previous text	Scotland is [summarized description...]
Continuation	One way to grow tomatoes is to	start with seeds...
Question answering	How many moons does Earth have?	Earth has one moon.
Chat	<i>Setup, followed by messages...</i>	<i>A sequence of relevant responses</i>

# Testing models in Azure OpenAI Studio playground

The screenshot shows the Azure OpenAI Studio interface, specifically the Chat playground section. The left sidebar includes links for Azure OpenAI, Playground, Chat (selected), Completions, DALL-E, Assistants (Preview), Management, Deployments, Models, Data files, Quotas, and Content filters (Preview). The main area has tabs for Replay chat, Clear chat, Playground Settings, View code, Show JSON, Import setup, Export setup, and Show panels. The Chat playground section features a 'Setup' panel with tabs for Prompt (selected) and Add your data, along with sections for Using templates, System message, and Examples. A central area contains a 'Start chatting' button and a text input field for user queries. To the right is a 'Configuration' panel with tabs for Deployment (selected) and Parameters, showing sliders for Max response (800), Temperature (0.7), Top P (0.95), Stop sequence, Frequency penalty (0), Presence penalty (0), and Current token count (11/4000). The bottom right corner shows a progress indicator for Input tokens.

## Exercise: Get started with Azure OpenAI Service



**Use the hosted lab environment if provided, or view the lab instructions at the link below:**

<https://aka.ms/mslearn-get-started-azure-openai>

# Knowledge check



**1** What do you need in order to test a generative AI model using the Azure OpenAI Service Studio?

- A deployed model name and Azure command line interface
  - An Azure OpenAI resource and an Azure Cognitive Services resource
  - An Azure OpenAI resource, a deployed model, and a playground
- 

**2** Which parameter could you adjust to change the randomness or creativeness of completions

- Temperature
  - Frequency penalty
  - Stop sequence
- 

**3** Which Azure OpenAI Studio playground is able to support conversational interactions that consist of a sequence of messages?

- Completions
- Chat
- Bot

# Use Azure OpenAI APIs in your applications



# Integrating Azure OpenAI into your app

Applications submit prompts to deployed models. Responses are completions.

- **Completion** - model takes an input prompt, and generates one or more predicted completions
- **Embeddings** - model takes input and returns a vector representation of that input
- **ChatCompletion** - model takes input in the form of a chat conversation (where roles are specified with the message they send), and the next chat completion is generated

**ChatCompletion** will be the endpoint we focus on for this course

Use **Completion** and **Embeddings** with GPT-3 based models

Use **ChatCompletion** with GPT-35-Turbo and later models

# Using the Azure OpenAI REST API

## Completion Endpoint

<https://endpoint.openai.azure.com/openai/deployments/deployment/completions>

```
{  
  "prompt": "Your favorite Shakespeare  
          play is",  
  "max_tokens": 5  
}
```



```
{  
  "id": "1234...",  
  "object": "text_completion",  
  "created": 1679001781,  
  "model": "gpt-35-turbo",  
  "choices": [  
    {  
      "text": "Macbeth",  
      "index": 0,  
      "logprobs": null,  
      "finish_reason": "stop"  
    }  
  ]  
}
```

# Using the Azure OpenAI REST API

## Embedding Endpoint

<https://endpoint.openai.azure.com/openai/deployments/deployment/embeddings>

```
{  
  "input": "The food was delicious and  
          the waiter was very  
          friendly..."  
}
```



```
{  
  "object": "list",  
  "data": [  
    {  
      "object": "embedding",  
      "embedding": [  
        0.0172990688066482523,  
        ....  
        0.0134544348834753042,  
      ],  
      "index": 0  
    }  
  ],  
  "model": "text-embedding-ada:002"  
}
```

# Using the Azure OpenAI REST API

## ChatCompletion Endpoint

<https://endpoint.openai.azure.com/openai/deployments/deployment/chat/completions>

```
{  
  "messages": [  
    {"role": "system",  
     "content": "You are an assistant  
       that teaches people about AI."},  
    {"role": "user",  
     "content": "Does Azure OpenAI  
       support multiple languages?"},  
    {"role": "assistant",  
     "content": "Yes, Azure OpenAI  
       supports several languages."},  
    {"role": "user",  
     "content": "Do other Cognitive  
       Services support translation?"}  
  ]  
}
```



```
{  
  "id": "unique_id", "object": "chat.completion",  
  "created": 1679001781, "model": "gpt-35-turbo",  
  "usage": { "prompt_tokens": 95,  
            "completion_tokens": 84, "total_tokens": 179},  
  "choices": [  
    {"message": { "role": "assistant",  
                "content": "Yes, other Azure Cognitive  
                  Services also support translation..."},  
    "finish_reason": "stop",  
    "index": 0}  
  ]  
}
```

# Using the Azure OpenAI SDKs

Language specific SDKs are available for use in your applications, in both C# and Python.

Code structure follows a similar pattern for both languages.

Parameters such as *Max Response* and *Temperature* are defined in the chat options.

Both synchronous and asynchronous API versions are available.

Pseudo code structure:

```
<include library>

<create client>

<define chat messages and options>

<send request>

<extract response content>
```

## Exercise: Integrate Azure OpenAI into your app



Use the hosted lab environment if provided, or view the lab instructions at the link below:

<https://aka.ms/mslearn-azure-openai-api>

## Knowledge check



1

Which REST endpoint should you use to interact with a GPT-4 model?

- Completion
- Embeddings
- ChatCompletion

2

When using the .NET SDK, which method should you use to call the ChatCompletion API?

- ChatMessage()
- GetChatCompletions()
- GetCompletions()

3

When using the Python SDK, which method should you use to call the ChatCompletion API?

- openai.ChatCompletion.get()
- openai.ChatCompletion.create()
- openai.chat.complete()

# Prompt engineering

---

Model completes the task by predicting the most probable next piece of text.

**Three main approaches:**

**Few-shot** – user includes several examples

**One-shot** – user provides one example

**Zero-shot** – no examples are provided

Convert the questions to a command:

Q: Ask Constance if we need some bread.

A: send-msg `find constance` Do we need some bread?

Q: Send a message to Greg to figure out if things are ready for Wednesday.

A: send-msg `find greg` Is everything ready for Wednesday?

Q: Ask Ilya if we're still having our meeting this evening.

A: send-msg `find ilya` Are we still having a meeting this evening?

Q: Contact the ski store and figure out if I can get my skis fixed before I leave on Saturday.

A: send-msg `find ski store` Would it be possible to get my skis fixed before I leave on Saturday?

Q: Thank Nicolas for lunch.

A: send-msg `find nicolas` Thank you for lunch!

Q: Tell Constance that I won't be home before 19:30 tonight – unmovable meeting.

A: send-msg `find constance` I won't be home before 19:30 tonight. I have a meeting at 19:30.

Q: Tell John that I need to book an appointment at 10:30.

A:

A few-shot prompt where multiple examples are provided. The model will generate the last answer

# Models

GPT-4

GPT-3.5

GPT3 base models

Embeddings series

DALL-E(Preview)

Whisper (Preview)

Fine-tuning models

Azure AI Studio > Models

## Models

Azure OpenAI is powered by models with different capabilities and price points. Deploy one of the provided base models for better performance and more accurate results.

[Learn more about the different types of base models](#)

Model name	Model version	Created at
babbage-002	1	9/5/2023 5:30 AM
davinci-002	1	9/5/2023 5:30 AM
gpt-35-turbo	0613	6/19/2023 5:30 AM
gpt-35-turbo-16k	0613	6/19/2023 5:30 AM
text-embedding-ada-002	2	4/3/2023 5:30 AM
whisper	001	9/14/2023 5:30 AM

Existing Azure OpenAI customers can apply using this form: <https://aka.ms/oai/get-gpt4>

# Prompts & Completions

Core component of the API service.

Text-in, text-out interface.

Provide an input **prompt**, and the model will generate a text **completion**.

The screenshot shows the Azure AI Studio interface with the 'Completions playground' selected in the sidebar. The main area displays five news headlines with their respective categories: Entertainment, Business, Tech, Politics, and Sports. The interface includes navigation links like 'Azure AI Studio' and 'Privacy & cookies', and buttons for 'Generate', 'Undo', 'Regenerate', and 'Tokens: 145'.

Azure AI Studio > Completions playground

Privacy & cookies

Completions playground

Deployments Examples

gpt-35-turbo Classify Text

View code

Classify the following news headline into 1 of the following categories: Business, Tech, Politics, Sport, Entertainment

Headline 1: Donna Steffensen Is Cooking Up a New Kind of Perfection. The Internet's most beloved cooking guru has a buzzy new book and a fresh new perspective  
Category: Entertainment

Headline 2: Major Retailer Announces Plans to Close Over 100 Stores  
Category: Business

Headline 3: The New York Times is Now Available on Snapchat Discover  
Category: Tech

Headline 4: Trump Administration to End DACA Program That Protected Young Undocumented Immigrants  
Category: Politics

Headline 5: The NFL's First Female Coach is Making History

Generate Undo Regenerate Tokens: 145

## Prompt types

Task type	Prompt example	Completion example
Classifying content	Tweet: I enjoyed the trip. Sentiment:	Positive
Generating new content	List ways of traveling	1. Bike 2. Car ...
Holding a conversation	A friendly AI assistant	
Transformation (translation and symbol conversion)	English: Hello French:	bonjour
Summarizing content	Provide a summary of the content {text}	The content shares methods of machine learning.
Picking up where you left off	One way to grow tomatoes	is to plant seeds.
Giving factual responses	How many moons does Earth have?	One

# Text tokens

Tokens

I have an orange cat named Butterscotch.

I have an orange cat named Butterscotch.

Horses are my favorite

animal	49.65%
animals	42.58%
\n	3.49%
!	0.91%

Probabilities:

IF TEMPERATURE IS 0

Horses are my favorite animal  
Horses are my favorite animal  
Horses are my favorite animal  
Horses are my favorite animal

IF TEMPERATURE IS 1

Horses are my favorite animal  
Horses are my favorite animals  
Horses are my favorite !  
Horses are my favorite animal

## Prompt Instruction

Suggest three names for an animal that is a superhero.

Animal: Cat

Names: Captain Sharpclaw, Agent Fluffball, The Incredible Feline

Animal: Dog

Names: Ruff the Protector, Wonder Canine, Sir Barks-a-Lot

Animal: Horse

Names:

## Completion Temperature 0 (always the same)

Mighty Equine, The Great Galloper, Thunderhoof

## Completion Temperature 1 (often different)

Blaze the Miracle Mare, Pegasus the Winged Warrior,  
Secretariat the Superhorse

## Completion Temperature 1 (often different)

Blaze of Glory, Sterling Silver, Thunderbolt

# Image tokens

## Cost depends on

- Size
- Detail setting
  - Low resolution mode
  - High resolution mode

### Low resolution mode

Fewer input tokens

Cost 85 tokens each

4096 x 8192 image – Cost is a fixed 85 tokens

### High resolution mode

Uses more tokens

Token cost is calculated by a series of scaling steps

2048 x 4096 image – Cost is 1105 tokens after calculation

# Azure OpenAI Service models

Models are grouped by family and capability

Model family	Description
<b>GPT-4</b>	A set of models that improve on GPT-3.5 and can understand generate natural language and code.
<b>GPT-3.5</b>	A set of models that improve on GPT-3 and can understand generate natural language and code.
<b>GPT3 base models</b>	Babbage-002 and Davinci-002 are GPT3 base models, intended for completion use cases.
<b>Embeddings</b>	A set of models that can convert text into numerical vector form to facilitate text similarity.
<b>DALL-E(Preview)</b>	A series of models that can generate original images from natural language.
<b>Whisper(Preview)</b>	A series of models in preview that can transcribe and translate speech to text.
<b>Fine-tuning models</b>	Let the customers tailor the Azure OpenAI models to the personal datasets.

# Azure OpenAI Models

# Models

<b>Text, code and Vision</b>	GPT-4 Turbo with Vision
<b>Text and code</b>	GPT-4 and GPT-4 Turbo Preview GPT-3.5, gpt-35-turbo and gpt-35-turbo-instruct, Babbage-002, Davinci-002
<b>Image</b>	DALL-E (Preview)
<b>Voice</b>	Whisper (Preview)
<b>Embeddings</b>	Embeddings
<b>Fine tuning</b>	GPT-3.5-Turbo (0613) Babbage-002 Davinci-002

## Multi modality

---

Multimodal AI systems can analyze, synthesize, and generate across text, images.



# Text and Code Models

# Introduction to Text and Code models



## GPT-3.5 models

Can understand and generate natural language and code



## GPT-4 models

Can solve difficult problems with greater accuracy

# GPT Model Use Cases

## GPT-3.5

Large, pretrained Language Models that use deep learning to generate content



## GPT-4

Built-in responsible AI to detect and mitigate harmful use



### Use cases:

- Generating natural language for chatbots and virtual assistants with awareness of the previous history of chat
- Power chatbots that can handle customer inquiries, provide assistance, and converse but doesn't have memory of conversations
- Automatically summarize lengthy texts
- Assist writers by suggesting synonyms, correcting grammar and spelling errors, and even generating entire sentences or paragraphs
- Help researchers by quickly processing large amounts of data and generating insights, summaries, and visualizations to aid in analysis
- Generate good quality code based on natural language

### Use cases:

- Generating and understanding natural language for customer service interactions, chatbots, and virtual assistants – doesn't have memory of conversations
- Generating high-quality code for programming languages based on natural language input.
- Providing accurate translations between languages
- Improving text summarization and content generation
- Provides for multi-modal interaction (text and images)
- Substantial reduction in Hallucinations
- Consistency between different runs is high

## Text and Code models- Key Terms

Term	Definition
▶ <b>Prompt</b>	<p>The text you send to the service in the API call. This text is then input into the model. For example, one might input the following prompt:</p>
	<p>Convert the questions to a command: Q: Ask Constance if we need some bread A: send-msg 'find constance' Do we need some bread?</p>
▶ <b>Completion or Generation</b>	<p>The text Azure OpenAI outputs in response. For example, the service may respond with the following answer to the above prompt: send-msg 'find greg' figure out if things are ready for Wednesday.</p>
▶ <b>Token</b>	<p>Azure OpenAI processes text by breaking it down into tokens. Tokens can be words or just chunks of characters. For example, the word hamburger gets broken up into the tokens ham, bur and ger, while a short and common word like pear is a single token. Many tokens start with a whitespace, for example hello and bye.</p>

# GPT-3.5 Turbo Instruct overview

Can understand and generate natural language or code.

This model is available to use with the Completions API.

The screenshot shows the Azure OpenAI service interface. On the left, a sidebar lists various features: Azure OpenAI, Playground, Chat, Completions, DALL-E (Preview), Management, Deployments, Models, Data files, Quotas, and Content filters (Preview). The main area is titled "Welcome to Azure OpenAI service" with the sub-headline "Explore the generative AI models, craft unique prompts for your use cases, and fine-tune". It displays a message "No deployment detected" with a "Create new deployment" button. Below this, there are two sections: "Get started" with "Text generation" icons and "Chat playground" (with a "Try it now" button) and "Completions playground" (with a "Try it now" button, highlighted with a red border).

# GPT3.5 overview

GPT-3.5 models can understand and generate natural language or code.

GPT-3.5 Turbo, has been optimized for chat and works well for traditional completions tasks as well.

Use the Chat Completions API to use GPT-3.5 Turbo.

Azure AI Studio

## Welcome to Azure OpenAI service

Explore the generative AI models, craft unique prompts for your use cases, and fine-tune

### Get started



#### Chat playground

Design a customized AI assistant using ChatGPT. Experiment with GPT-3.5-Turbo and GPT-4 models.

[Try it now](#)



#### Completions playground

Experiment with completions models for use cases such as summarization, content generation, and classification.

[Try it now](#)

## GPT-3.5 models

Model ID	Max Request (tokens)
gpt-35-turbo(0301)	4,096
gpt-35-turbo (0613)	4,096
gpt-35-turbo-16k (0613)	16,384
gpt-35-turbo-instruct (0914)	4097
gpt-35-turbo (1106)	Input: 16,385 Output: 4,096

## GPT-4 Overview

A set of models that improve on GPT-3.5 and can understand as well as generate natural language and code.

A large multimodal model (accepting text inputs, emitting text outputs) that,

Exhibits human-level performance on various professional and academic benchmarks.

GPT-4 is more reliable, creative, and able to handle much more nuanced instructions than GPT-3.5.

Can do everything that GPT 3.5 can do

Visual Question and Answering (VQA)

Steerability

context length of 8,192 tokens

Apply here for GPT-4 access(for already onboarded customers): <https://aka.ms/oai/get-gpt4>

## GPT-4 models

Model ID	Max Request (tokens)
gpt-4 (0314)	8,192
gpt-4-32k (0314)	32,768
gpt-4 (0613)	8,192
gpt-4-32k (0613)	32,768
gpt-4 (1106-preview) GPT-4 Turbo Preview	Input: 128,000 Output: 4096
gpt-4 (vision-preview) GPT-4 Turbo with Vision Preview	Input: 128,000 Output: 4096

## Babbage-002 and Davinci-002

Intended for completion  
use cases

Generate natural language  
or code

Not trained for instruction  
following

# Working with the Chat models

There are two different APIs available for working with these type of models.

## Chat Completion API (Recommended)

A new API designed for working with Chat models (gpt-35-turbo and gpt-4). We recommend using this API moving forward.

In this API, you pass in a series of messages rather than the prompt as a raw string making it easier to manage the conversation.

```
messages = [  
    {  
        "role": "system",  
        "content": "Assistant is a large language model trained by OpenAI."  
    },  
    {  
        "role": "user",  
        "content": "What is a garbanzo bean?"  
    }  
]
```

## Completion API

The GPT-35-Turbo model also works with the existing completions API. In this API, you'll format your prompt as a string and need to use the special tokens from [Chat Markup Language](#).

We generally recommend using the Chat API instead because the special tokens may get updated in future version of gpt-35-turbo meaning you would need to update your prompts before updating to newer versions of the model.

```
<|im_start|>system  
Assistant is a large language model trained by OpenAI.  
<|im_end|>  
<|im_start|>user  
What is a garbanzo bean?  
<|im_end|>  
<|im_start|>assistant
```

# Working with the Chat models

## Previous GPT-3 models

Previous models were text-in and text-out

(i.e., they accepted a prompt string and returned a completion to append to the prompt).

---

Answer questions from the context below.

Context:

A neutron star is the collapsed core of a massive supergiant star, which had a total mass of between 10 and 25 solar masses, possibly more if the star was especially metal-rich.

Q: What is a neutron star?

A:

## GPT-35-Turbo and GPT-4

The GPT-35-Turbo and GPT-4 models are conversation-in and message-out.

---

(i.e., it expects a prompt string that is formatted in a specific chat-like transcript format and returns a completion that represents a model-written message in the chat)

---

### [System Message]

Assistant is an AI Chatbot designed to answer questions from the context provided below.

Context:

A neutron star is the collapsed core of a massive supergiant star, which had a total mass of between 10 and 25 solar masses, possibly more if the star was especially metal-rich.

### [User]

What is a neutron star?

### [Assistant]

# An example in Chat playground

The screenshot shows the Azure AI Studio interface with the "Chat playground" selected in the sidebar. The main area displays the "Assistant setup" and "Chat session" components.

**Assistant setup:** This section includes a "System message" input field and a "Add your data (preview)" button. It also contains a note about securely storing data in the Azure subscription and a link to learn more about data protection. A "Add a data source" button is present at the bottom.

**Chat session:** This section shows a user message "What is a neutron star?" and an AI response describing a neutron star as a dense astronomical object formed from collapsed stars. It also includes "Clear chat", "View code", and "Show raw JSON" buttons.

**User message:** A text input field at the bottom with placeholder text "Type user query here. (Shift + Enter for new line)".

# Understanding the prompt format

## The system message

The system message is included at the beginning of the prompt and is used to prime the model and you can include a variety of information in the system message including:

- A brief description of the assistant
- The personality of the assistant
- Instructions for the assistant
- Data or information needed for the model

## User and assistant messages

After the system message, you can include a series of messages between the *user* and the *assistant*. You denote who the message is from by setting the role to user or assistant.

```
{  
  "role": "user",  
  "content": "What is a garbanzo bean?"  
}
```

## Example prompt

```
{ "role": "system", "content": "You are an Xbox customer support agent whose primary goal is to help users with issues they are experiencing with their Xbox devices. You are friendly and concise. You only provide factual answers to queries, and do not provide answers that are not related to Xbox. },  
{ "role": "user", "content": "Why won't my Xbox turn on?"},  
{ "role": "assistant", "content": "There could be a few reasons why your Xbox isn't turning on...."},  
{ "role": "user", "content": "I confirmed the power cord is plugged in but it's still not working" }
```

# GPT-4 Turbo with Vision

# GPT-4 Turbo with Vision

A multimodal model that analyze images and provide textual responses to questions about them

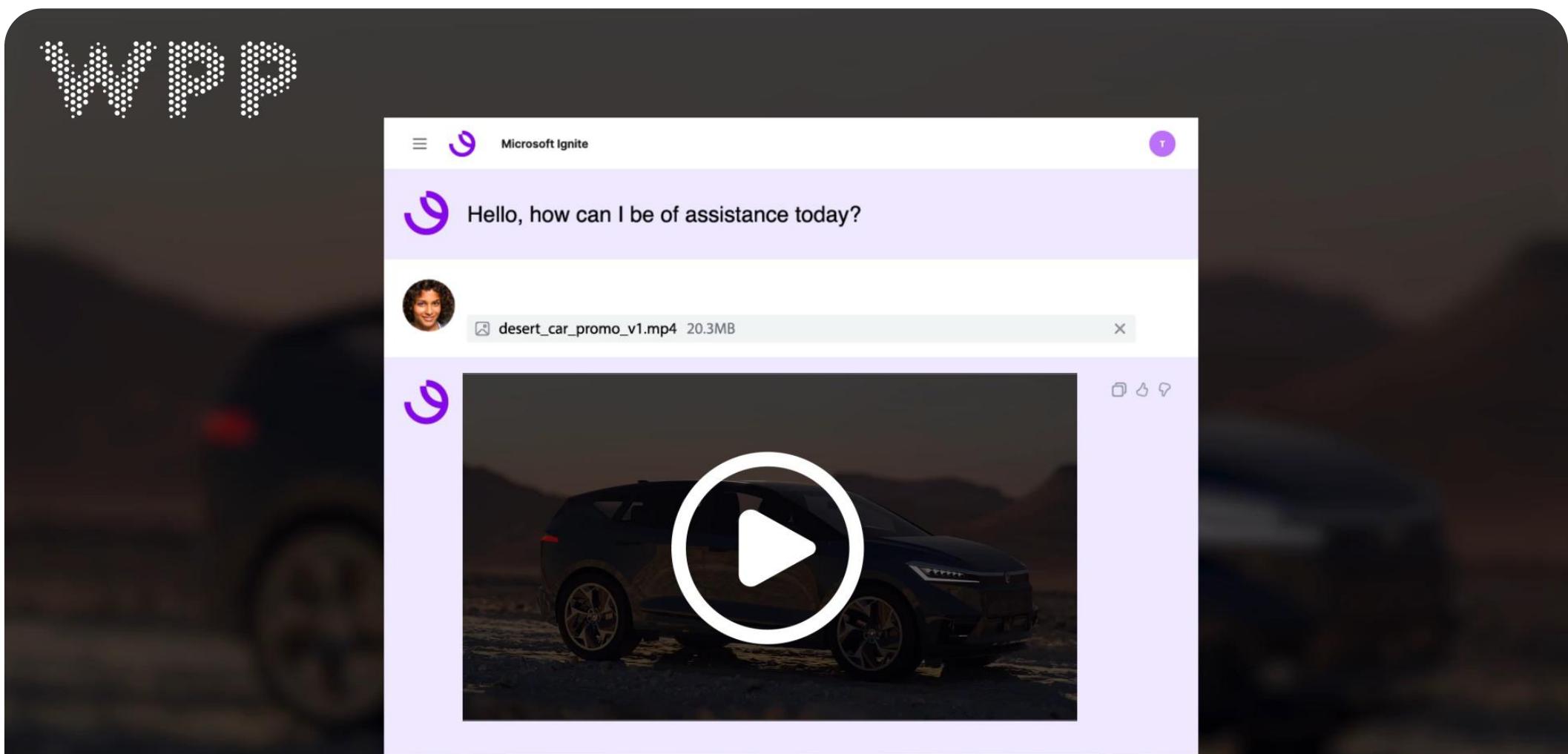
Incorporates both natural language processing and visual understanding.

The screenshot shows the Azure AI Studio interface with the 'Chat playground' selected. On the left, a sidebar lists 'Azure OpenAI' and 'Playground' sections, with 'Chat' highlighted. The main area is titled 'Chat playground' and contains an 'Assistant setup' dialog and a 'Chat session' panel.

**Assistant setup:** This dialog is titled 'Assistant setup' and has a tab for 'System message'. It includes a 'Save changes' button and a section for specifying how the chat should act. It suggests using a template or writing a system message. A red box highlights the 'System message' input field, which contains the text: "You are an AI assistant that helps people find information."

**Chat session:** This panel is titled 'Chat session' and includes buttons for 'Clear chat', 'Playground Settings', 'View code', and a toggle for 'raw JSON'. It features a 'Start chatting' section with a robot icon and instructions for using the microphone and speaker icons. Below it is a text input field labeled 'Describe this image!' with three red-bordered buttons at the bottom right.

# GPT-4 Turbo with Vision + Azure AI Services





## Tents

UNFOLD ADVENTURE WITH CONTOSO'S RANGE OF TENTS - YOUR PORTABLE HOME UNDER THE stars. From intimate two-person camping shelters to family-sized canopies, our tents are your ticket to escaping the ordinary and embracing the great outdoors.



TrailMaster X4 Tent



Alpine Explorer Tent



SkyView 2-Person Tent



# Combine Text & Image Inputs

## User

Write a six-line rhyming poem. Starting each line with the letter "c".

The first three lines are around this image.



And the next three lines are around this image.



## Vision AI

Calm and cool, the sea is blue,  
Crashing waves, a scenic view,  
Craggy rocks, a sight so true,  
Colorful and warm, the sun is low,  
Casting shadows, a golden glow,  
Cozy and calm, a peaceful show.

# Detailed Image Descriptions



## Vision AI

A person with a striped shirt and an apron is holding a metal watering can and watering the plants in the greenhouse.

A large metal watering can with a long spout is being used by a person to water the green and leafy plants in the greenhouse.

A greenhouse full of green and leafy plants is being watered by a person with a metal watering can. The person is wearing a striped shirt and an apron.

A person is taking care of the plants in the greenhouse by watering them with a metal watering can. The person has a striped shirt and an apron on. The plants are green and leafy.

# DALL-E 3(Preview)

# DALL-E (Preview)

The model that works with images.  
Image capabilities include,

**Image Generation**

**Editing an image**

**Creating variations of image**

Dall-E 3 is now available in preview

The screenshot shows the DALL-E playground (Preview) interface within the Azure AI Studio. The left sidebar lists various features: Azure OpenAI, Playground, Chat, Completions, DALL-E (Preview) (which is selected and highlighted in grey), Management, Deployments, Models, Data files, Quotas, and Content filters (Preview). The main area has a header "DALL-E playground (Preview)" with back and forward navigation arrows, and links to "View code" and "Settings". Below the header is a "Prompt" section with the text "Watercolour painting of Seattle skyline". Two generated images are displayed below: a watercolor painting of the Seattle skyline at sunset with a yellow-orange sky and a more vibrant, colorful version of the same scene. Each image has a caption "Watercolour painting of Seattle skyline" and a row of five small icons for download, copy, save, delete, and refresh.

# DALL-E 3

*Azure OpenAI Service*

DALL-E 3 is an image generation model that allows you to generate images from text prompts



# DALL·E 3 brings dramatic quality improvements

Enhanced image quality



*"an image of a bear walking through the forest, close-up, photorealistic"*

Ability to render more complex scenes



*"An astronaut lounging on a chair made of clouds with the vivid night sky and earth brightly shining in the background. A large meteor streaks through the sky as the astronaut watches it fly by."*

More realistic images of people and hands



*"A group of business people meet in a conference room of a high-rise with a stunning city-scape in the background"*

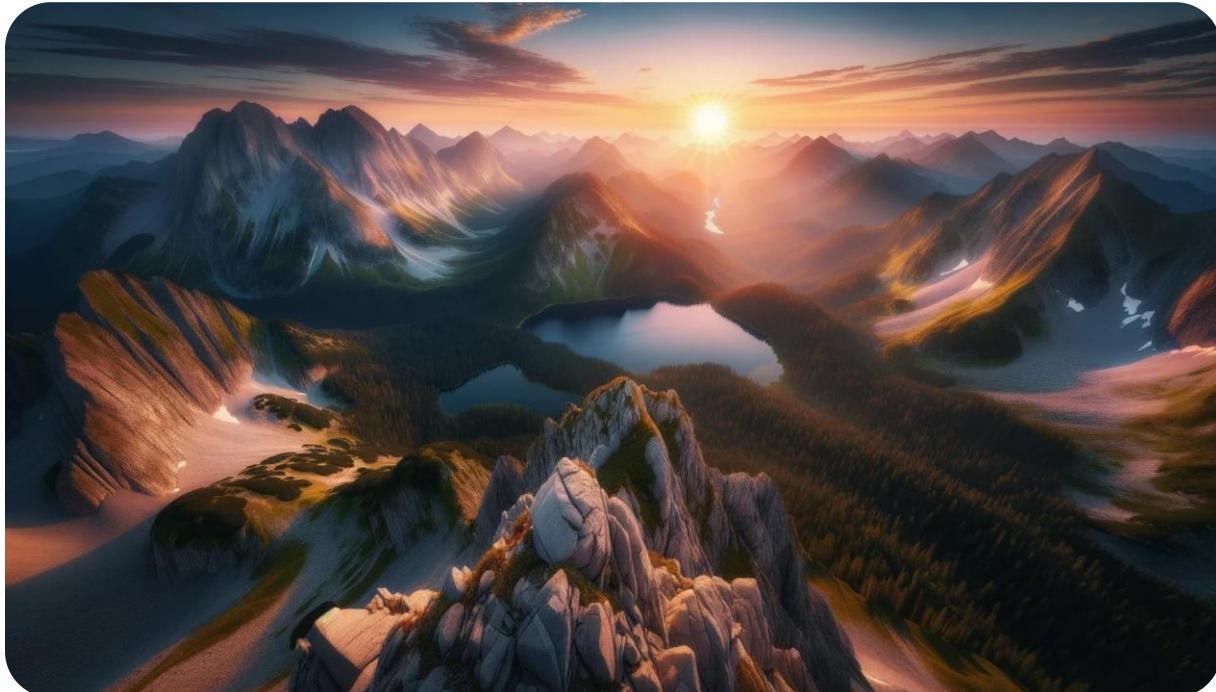
# DALL·E 3 also offers new capabilities

Improved rendering of text in images



*"a 3d render of a cute robot telling weary travelers "I hope you find what you're looking for", futuristic background"*

Support for new aspect ratios



*"A panorama of vast mountains with vivid mountain peaks, forests, and alpine lakes at sunrise"*

# Prompt rewriting in DALL·E 3

The DALL-E 3 API will include built-in prompt rewriting to enhance images, reduce bias, and increase natural variation



*"A group of scientists in the lab"*



*"A painting of mount rainier, various styles"*



# Use Cases for DALL·E 3



**LOGO & BRANDING:**  
QUICK CONCEPT  
GENERATION.



**CREATIVE  
INSPIRATION:**  
OVERCOME DESIGN  
BLOCKS.



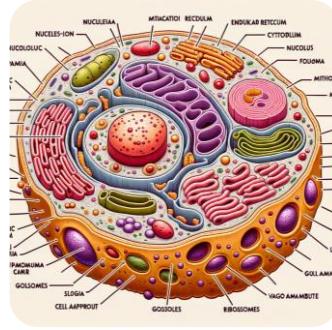
**CONTENT  
ILLUSTRATIONS:**  
UNIQUE IMAGES FOR  
BLOGS/ARTICLES.



**AD CAMPAIGNS:**  
VISUALIZE MARKETING  
CONCEPTS.



**PRODUCT  
VISUALIZATION:** GAUGE  
INTEREST & FEEDBACK.



**EDUCATION:** CUSTOM  
IMAGERY FOR COURSES.



**FASHION DESIGN:**  
VISUALIZE CLOTHING  
PATTERNS.



**GAMING:** CHARACTER &  
ENVIRONMENT  
CONCEPTS.

# Content Credentials

Provides information about the origin of an image generated by the DALL-E series models.

- Ways to check the Credential,
  - Content Credentials Verify webpage
  - Content Authenticity Initiative (CAI) JavaScript SDK

The screenshot shows the Microsoft Content Credentials interface. On the left, there's a file upload area with a placeholder: "Select another file from your device or drag and drop anywhere". Below it is a preview of an image titled "Untitled asset" from Nov 6, 2023. A "Search for possible matches" button is nearby. At the bottom left is a "Change language" dropdown. On the right, there's a detailed view of an image of a white animal in a snowy landscape. This view includes the title "Untitled asset", the date "Nov 6, 2023", and a "Compare" button. To the right of the image are sections for "Process" (with a dropdown arrow), "About this Content Credential" (with a dropdown arrow), "Issued by" (listing "Microsoft Corporation" with a question mark icon), and a note stating "This is the trusted organization, device, or individual that recorded the details above and issued this Content Credential." At the bottom right is an "Issued on" section with the date "Nov 6, 2023 at 2:25 PM PST".

# Provisioned throughput units(PTU)

# What is Provisioned Throughput?



A new Azure OpenAI Service feature that lets customers **reserve model processing capacity** for running high-volume or latency-sensitive workloads



Reserved processing capacity provides consistent latency and throughput for workloads with consistent characteristics, such as prompt size, completion size, and number of concurrent API requests



Processing capacity is defined in units called "Provisioned Throughput Units" (PTUs) that are purchased on a monthly commitment



Once purchased, customers use PTUs to create provisioned Azure OpenAI deployments of GPT-35-turbo or GPT-4 models during the term of their commitment

## Provisioned Throughput Key Values

Model processing capacity for your high-volume production workload



### Predictable performance

Stable latency  
and throughput for  
uniform workloads



### Reserved Processing Capacity

Processing capacity is available  
to meet demand



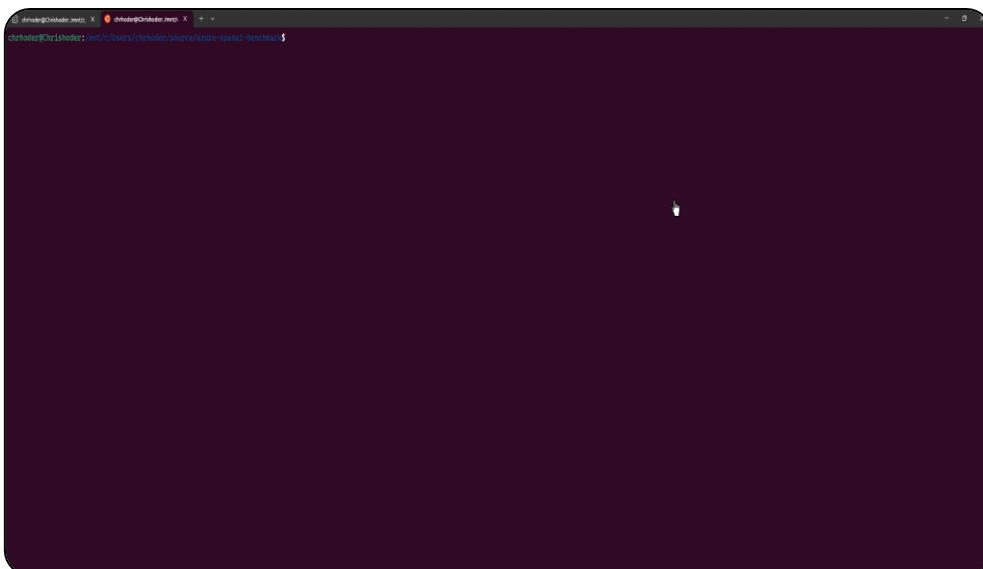
### Cost Savings

High throughput workloads will  
likely see cost savings vs token-  
based consumption

# Managing Provisioned Deployments

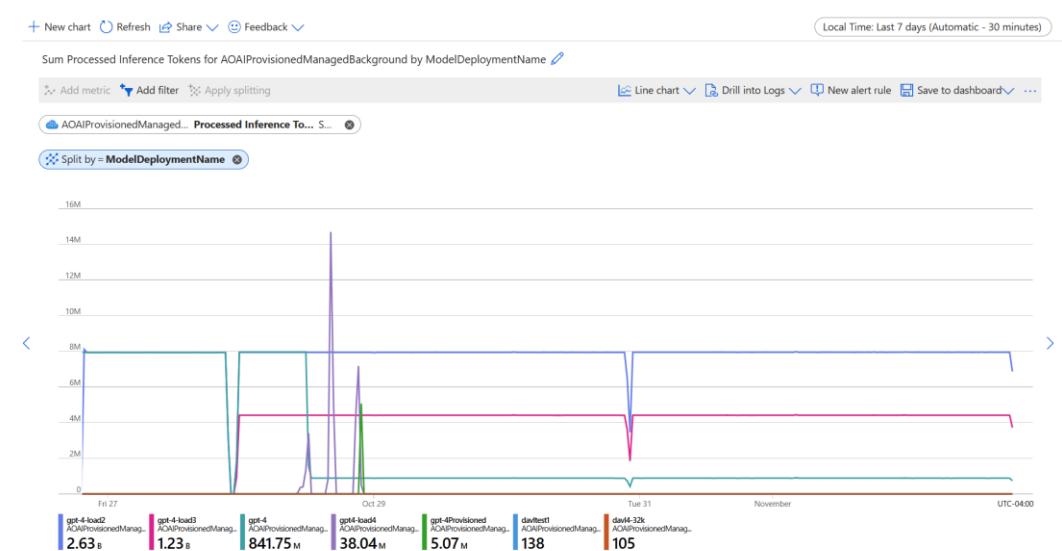
## Performance Assessment

Azure OpenAI provides a Python toolset for assessing deployment performance across the range of prompt/generation token sizes and RPM

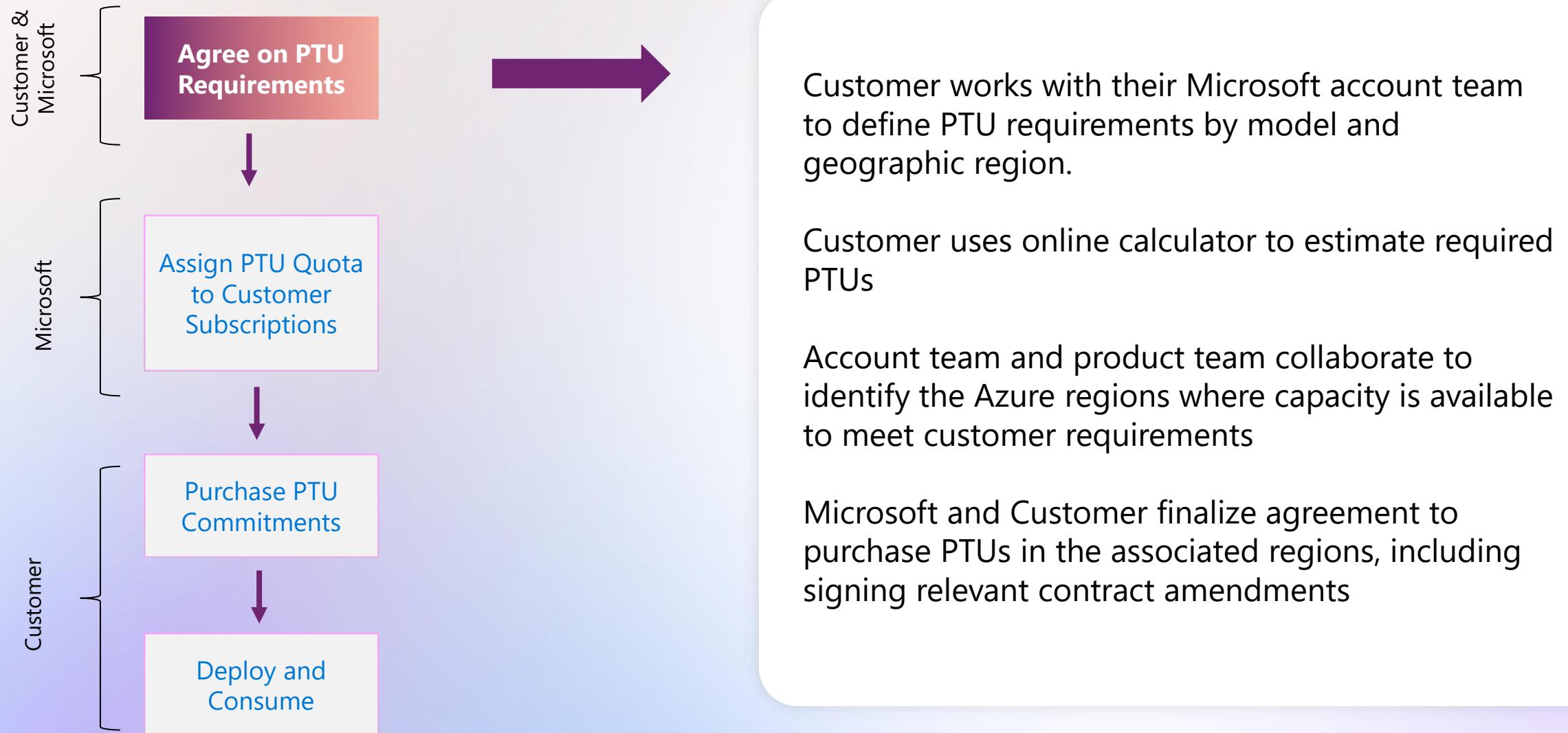


## Production Monitoring

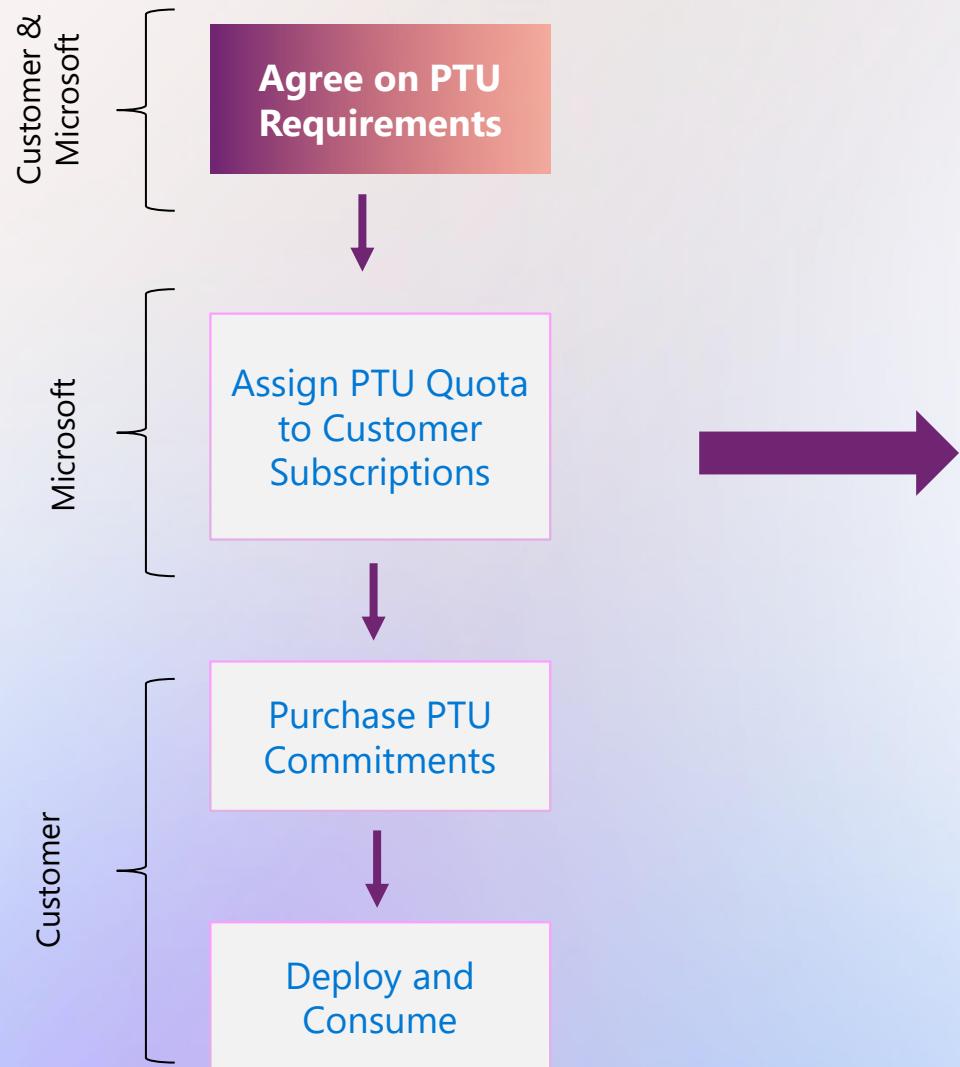
Deployment metrics are built into Azure Monitor, including utilization, latency, token and request counts



# Provisioned Throughput Workflow



# Provisioned Throughput Workflow

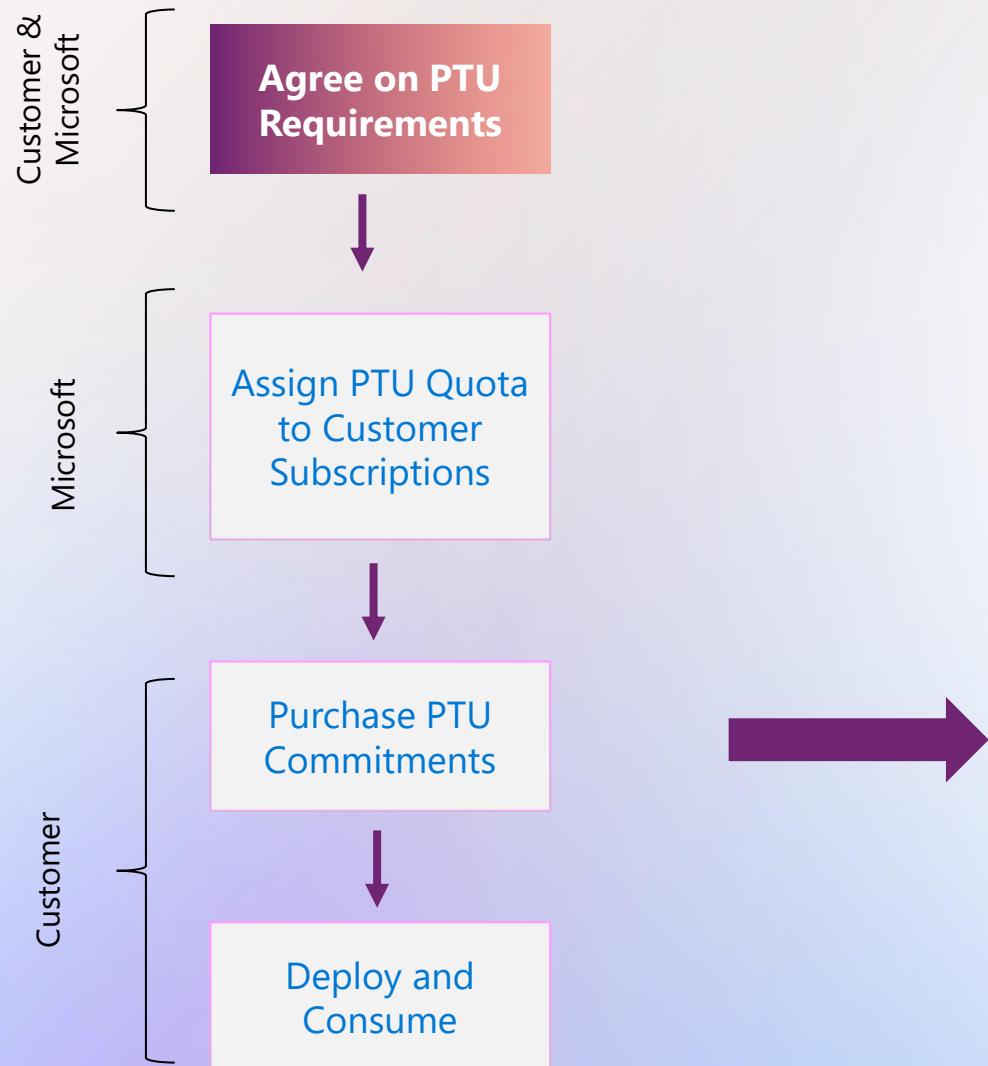


Microsoft product team reserves the capacity for the customer and assigns PTU quota to the subscription

PTU quota is assigned per-model type

The expectation is that the PTU quota will be purchased within 48 hours of assignment. Otherwise, PTUs may be reclaimed.

# Provisioned Throughput Workflow



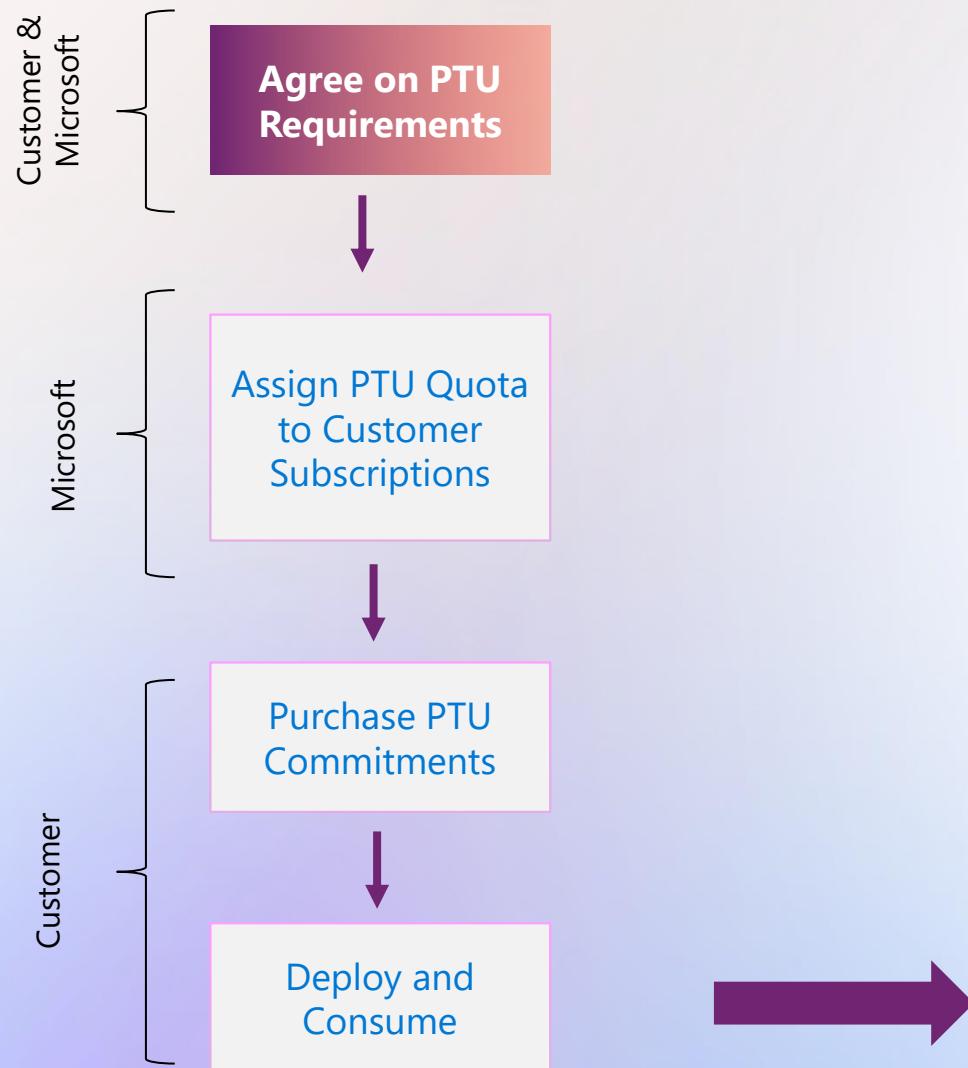
Customer uses Azure AI Studio to purchase the PTUs

Commitments are created at the resource level, not the subscription level, and cannot be transferred to other resources or subscriptions once created

Each resource can have its own commitment for a different number of PTUs (constrained by the sub quota limit)

Provisioned throughput deployments can only be created within resources that have commitments, and only up to the commitment amount

# Provisioned Throughput Workflow



Provisioned deployments are created in the same way as standard, pay-as-you-go deployments

The difference is that the deployment type will be set to "Provisioned Managed", and the desired number of PTUs will need to be set

Once a provisioned deployment is created, it can be scaled up or down in PTU increments.

Multiple provisioned deployments of the same or different model types can be created within the same resource

Assess performance using benchmarking tool

Monitor utilization via Azure Monitor

# Using function calling with Azure OpenAI Service (Preview)

# How to use function calling with Azure OpenAI Service (Preview)

Call the chat completions API with your functions and the user's input

Use the model's response to call your API or function

Call the chat completions API again, including the response from your function to get a final response

```
def run_conversation():
    # Step 1: send the conversation and available functions to the model
    messages = [{"role": "user", "content": "What's the weather like in San Francisco,"}
    tools = [
        {
            "type": "function",
            "function": {
                "name": "get_current_weather",
                "description": "Get the current weather in a given location",
                "parameters": {
                    "type": "object",
                    "properties": {
                        "location": {
                            "type": "string",
                            "description": "The city and state, e.g. San Francisco, CA"
                        },
                        "unit": {"type": "string", "enum": ["celsius", "fahrenheit"]},
                    },
                    "required": ["location"]
                }
            }
        }
    ]
```

# Parallel function calling

Allows to perform multiple function calls together

```
tool_calls = response_message.tool_calls
# Step 2: check if the model wanted to call a function
if tool_calls:
    # Step 3: call the function
    # Note: the JSON response may not always be valid; be sure to handle errors
    available_functions = {
        "get_current_weather": get_current_weather,
    } # only one function in this example, but you can have multiple
    messages.append(response_message) # extend conversation with assistant's response
    # Step 4: send the info for each function call and function response to the model
    for tool_call in tool_calls:
        function_name = tool_call.function.name
        function_to_call = available_functions[function_name]
        function_args = json.loads(tool_call.function.arguments)
        function_response = function_to_call(
            location=function_args.get("location"),
            unit=function_args.get("unit"),
        )
        messages.append(
            {
                "tool_call_id": tool_call.id,
                "role": "tool",
                "name": function_name,
                "content": function_response,
            }
        ) # extend conversation with function response
    second_response = client.chat.completions.create(
        model=<REPLACE_WITH_YOUR_1106_MODEL_DEPLOYMENT_NAME>,
        messages=messages,
    ) # get a new response from the model where it can see the function responses
    return second_response
print(run_conversation())
```

# Working with function calling

Defining functions

Managing the flow with functions

Prompt engineering with functions

```
response = openai.ChatCompletion.create(  
    deployment_id="gpt-35-turbo-0613",  
    messages=messages,  
    functions=functions,  
    function_call="auto",  
)  
response_message = response["choices"][0]["message"]  
  
# Check if the model wants to call a function  
if response_message.get("function_call"):  
  
    # Call the function. The JSON response may not always be valid so make sure to handle it  
    function_name = response_message["function_call"]["name"]  
  
    available_functions = {  
        "search_hotels": search_hotels,  
    }  
    function_to_call = available_functions[function_name]  
  
    function_args = json.loads(response_message["function_call"]["arguments"])  
    function_response = function_to_call(**function_args)  
  
    # Add the assistant response and function response to the messages  
    messages.append( # adding assistant response to messages  
    {  
        "role": response_message["role"],  
        "function_call": {  
            "name": function_name,  
            "arguments": response_message["function_call"]["arguments"],  
        },  
    },  
)
```

# Using function calling responsibly

Validate Function Calls

Use Trusted Data and Tools

Follow the Principle of Least Privilege

Consider Real-World Impact

Implement User Confirmation Steps

# Availability and pricing

# Region availability

Regions currently available at launch:

- Australia East
- Canada East
- East US
- East US 2
- France Central
- Japan East
- North Central US
- South Central US
- Sweden Central
- Switzerland North
- UK South
- West Europe

Products	EUROPE		FRANCE		UNITED STATES					
	North Europe	West Europe	France Central	France South	Central US	East US	East US 2	North Central US	South Central US	West Central US
Azure AI services	✓	✓	✓		✓	✓	✓	✓	✓	✓
Azure AI Anomaly Detector	✓	✓	✓		✓	✓	✓	✓	✓	✓
Azure AI Vision	✓	✓	✓		✓	✓	✓	✓	✓	✓
Content Moderator	✓	●	✓		✓	✓	✓	✓	✓	✓
Azure AI Custom Vision	✓	✓			●	✓	✓	✓	✓	✓
Face API	✓	✓	✓		✓	✓	✓		✓	✓
Azure AI Language	✓	✓	✓		✓	✓	✓	✓	✓	✓
Language Understanding (LUIS)	✓	✓	✓		✓	✓	✓	✓	✓	✓
Azure AI Personalizer	✓	✓	✓		✓	✓	✓	✓	✓	□
QnA Maker	✓	✓	✓		✓	✓	✓	✓	✓	✓
Speaker recognition										
Speech Services	✓	✓	✓		✓	✓	✓	✓	✓	✓
Translator										
Azure OpenAI Service		✓	✓			✓			✓	
Azure AI Immersive Reader	✓	✓	✓		✓	✓	✓	✓	✓	✓
Azure AI Metrics Advisor	✓	✓	✓		✓	✓	✓	✓	✓	
Azure AI Video Indexer	✓	✓	✓		✓	✓	✓	✓	✓	✓

\*Availability as of December 2023

## Pricing



Pricing will be based on the pay-as-you-go consumption model



Price per unit for each model



Similar to other Azure Cognitive Services pricing models

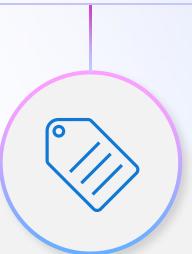
## General best practices to mitigate throttling during autoscaling



Implement retry logic in your application.



Avoid sharp changes in the workload.



Test different load increase patterns. Create another OpenAI service resource and distribute the workload among them.

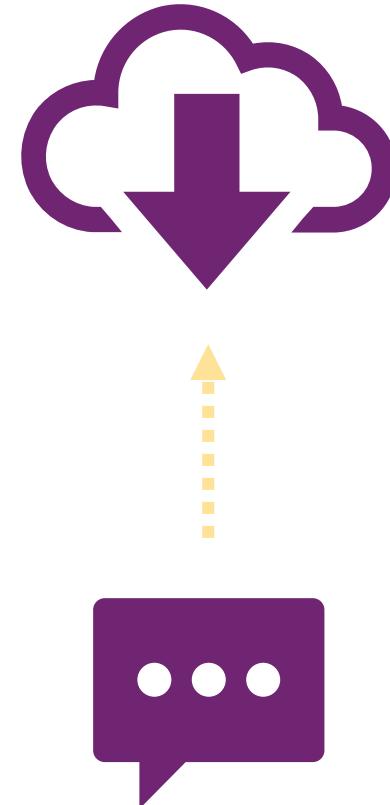
# Apply prompt engineering with Azure OpenAI Service



# What is Prompt Engineering?

Constructing prompts to:

- Maximize relevancy and accuracy of completions
- Specify formatting and style of completions
- Provide conversational context
- Mitigate bias and improve fairness



## Providing clear instructions

Write a product description for a new water bottle



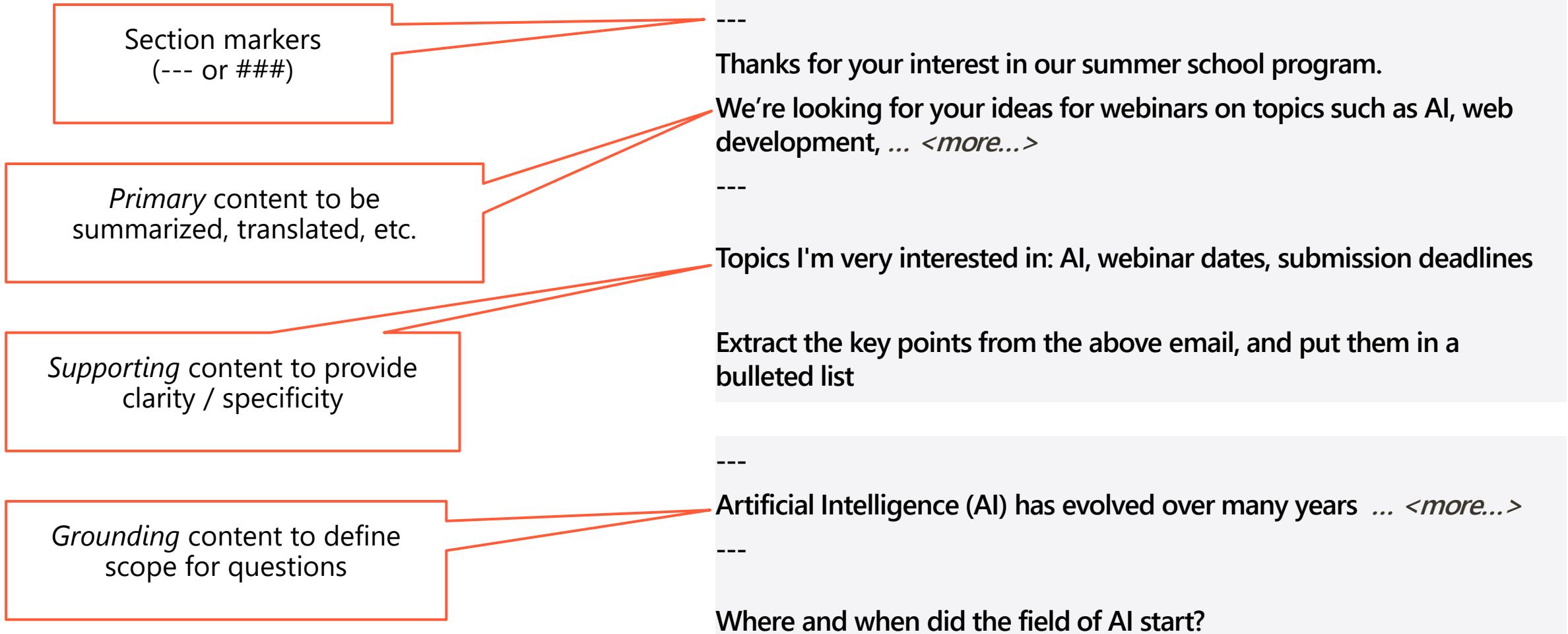
Introducing the latest addition to our product line - the innovative and eco-friendly water bottle. Made from high-quality, BPA-free materials, this bottle is designed to keep your beverages at the perfect temperature for hours. Its sleek and stylish design makes it the perfect accessory for any outdoor activity or workout session. With its convenient carrying handle and leak-proof lid, you can take it with you wherever you go, without worrying about spills or leaks. Plus, its easy-to-clean design ensures that you can use it again and again, making it the perfect choice for anyone looking for a reliable and sustainable water bottle.

Write a product description for a new water bottle that is 100% recycled. Be sure to include that it comes in natural colors with no dyes, and each purchase removes 10 pounds of plastic from our oceans



Introducing our newest water bottle, the eco-friendly choice for the conscious consumer. Made from 100% recycled materials and available in natural colors with no dyes, this bottle is not only stylish but also sustainable. Each purchase of this bottle removes 10 pounds of plastic from our oceans, making it a responsible choice for those who care about the environment. Its durable and leak-proof design makes it perfect for outdoor activities, while its easy-to-carry handle makes it convenient to take on the go. Join us in our mission to reduce plastic waste and choose this eco-friendly water bottle for your daily hydration needs.

## *Primary, supporting, and grounding content*



## Cues

- Cues provide a starting point on which the completion builds
- Use cues to set the scope and format for the output
- Cues are particularly useful when generating code.  
For example:
  - Provide the first keyword in a SQL SELECT statement
  - Define a name and parameters for a Python function

---

**Review: The writing was good but the movie was too long.**

**Review: The special effects were terrible.**

*<more>*...

---

**Summarize the reviews above.**

**Most common complaints:**

-



- **The movie was too long.**
- **The special effects were terrible.**

# Requesting output composition

Write a table in markdown with 6 animals in it, with their genus and species



Animal	Genus	Species
Red Panda	Ailurus	fulgens
African Elephant	Loxodonta	africana
Snow Leopard	Panthera	uncia
Arctic Fox	Vulpes	lagopus
Chimpanzee	Pan	troglodytes
Siberian Tiger	Panthera	tigris altaica

## Using a system message

When using the ChatCompletions API, a system message sets context that can affect completions

```
{"role": "system", "content": "You are a casual, helpful assistant.  
You will talk like an American old western film character."},  
{"role": "user", "content": "Can you direct me to the library?"}
```



```
{"role": "assistant", "content": "Well howdy there, stranger! The library, huh?  
Y'all just head down the main road till you hit the town  
square. Once you're there, take a left and follow the street  
for a spell. You'll see the library on your right, can't  
miss it. Happy trails!"}
```

## Conversation history and *few shot* learning

- Including previous messages helps retain the context and style of a conversation
- Starting a conversation with user-defined messages and responses (*few shot* learning) sets examples for the expected responses

```
{"role": "system", "content": "You are an assistant that evaluates the sentiment of customer feedback."},  
{"role": "user", "content": "That was an awesome experience"},  
{"role": "assistant", "content": "positive"},  
{"role": "user", "content": "I won't do that again"},  
{"role": "assistant", "content": "negative"},  
{"role": "user", "content": "That was not worth my time"},  
{"role": "assistant", "content": "negative"},  
{"role": "user", "content": "You can't miss this"}
```



```
{"role": "assistant", "content": "positive"},
```

# Chain of thought

Ask a model to break down its response and explain its reasoning

**What sport is easiest to learn but hardest to master? Give a step by step approach of your thoughts, ending in your answer**

Step 1: Identify the criteria for "easy to learn" and "hard to master" sports.

For a sport to be considered easy to learn, it should have simple rules and require minimal equipment. For a sport to be considered hard to master, it should require years of practice to perfect and have a large variety of techniques and strategies...

Step 2: Consider different sports that fit these criteria.

Some sports that could be considered easy to learn but hard to master include:

- Tennis: Simple rules, easy to pick up, but requires years of practice to master all the different shots and strategies.
- Golf: Basic swing mechanics are easy to learn but perfecting them takes a lifetime of practice.
- ...

Step 3: Evaluate each sport based on the criteria.

Tennis and golf both require expensive equipment, which could make them less accessible to beginners. In terms of skill ceiling and time required to master, tennis and golf both have a high skill ceiling and require years of practice to perfect....

Step 4: Make a decision based on the evaluation.

Based on the above criteria and evaluation, I would say that tennis is the sport that is easiest to learn but hardest to master. While it may require expensive equipment, it is still accessible to many people. The basic rules and techniques are easy to learn, but mastering all the different shots and strategies takes years of practice...

## Exercise: Utilize prompt engineering in your app



**Use the hosted lab environment if provided, or view the lab instructions at the link below:**

<https://aka.ms/openai-prompt-engineering-lab>

# Knowledge check



- 1 How can developers optimize the performance of Azure OpenAI models?**
  - By using complex instructions that are difficult to understand
  - By providing clear and descriptive instructions
  - By using vague prompts

---
- 2 What is the purpose of the system message in a prompt?**
  - To give the model instructions, perspective, or other information helpful to guide its response
  - To give the model a specific answer to generate
  - To provide filler information to the model

---
- 3 What is the purpose of providing conversation history to an AI model?**
  - Providing conversation history to an AI model is irrelevant and has no effect on the AI's performance.
  - To limit the number of input tokens used by the model
  - To enable the model to continue responding in a similar way and allow the user to reference previous content in subsequent queries

# Generate code with Azure OpenAI Service



## Use natural language prompts to generate code

- Use prompt engineering to describe what you want clearly
- Break down complex tasks into smaller pieces

### Prompt:

write a function for binary search in python

```
def binary_search(array, target):  
    low = 0  
    high = len(array) - 1  
  
    while low <= high:  
        mid = (low + high) // 2  
        if array[mid] == target:  
            return mid  
        elif array[mid] < target:  
            low = mid + 1  
        else:  
            high = mid - 1  
  
    return -1
```

# Change programming language and understand unknown code

## Change programming language

Convert the following code to C#:

```
def print_squares(n):
    for i in range(1, n+1):
        print(i**2)
```



```
public void PrintSquares(int n)
{
    for (int i = 1; i <= n; i++)
    {
        Console.WriteLine(i * i);
    }
}
```

## Understand unknown code

Could you explain what this code is doing?:

```
fn qzplv(n: i32, m: i32) -> i32 {
    if n == 0 {
        return m + 1;
    } else if m == 0 {
        return qzplv(n - 1, 1);
    } else {
        return qzplv(n - 1, qzplv(n, m - 1));
    }
}
```



Sure, I'd be happy to explain what this code is doing!

This code is implementing the Ackermann function,...

## Complete partial code

- Writing your prompt with the most amount of information available will be the most accurate
- If enough code is given to determine the programming language, the model will use that language

### Prompt:

complete the following function

"""

```
# calculate the average of the numbers
# in an array, but only if they're even
def
```

```
# calculate the average of the even
# numbers in an array
def calculate_even_average(numbers):
    total = 0
    count = 0
    for number in numbers:
        if number % 2 == 0:
            total += number
            count += 1
    if count == 0:
        return 0 # avoid division by zero
    else:
        return total / count
```

## Write unit tests

### Prompt:

write three unit tests for this function

"""

```
def binary_search(arr, target):
    left, right = 0, len(arr) - 1
    while left <= right:
        mid = (left + right) // 2
        if arr[mid] == target:
            return mid
        elif arr[mid] < target:
            left = mid + 1
        else:
            right = mid - 1
    return -1
```

```
def test_binary_search():
    # Test case 1: Target element is in the middle of the
    # array
    arr1 = [1, 2, 3, 4, 5, 6]
    target1 = 4
    assert binary_search(arr1, target1) == 3

    # Test case 2: Target element is at the beginning of
    # the array
    arr2 = [1, 2, 3, 4, 5, 6]
    target2 = 1
    assert binary_search(arr2, target2) == 0

    # Test case 3: Target element is not in the array
    arr3 = [1, 2, 3, 4, 5, 6]
    target3 = 7
    assert binary_search(arr3, target3) == -1
```

## Fix bugs in your code

- Help identify and fix bugs quickly
- Provide an explanation of the issues

### Prompt:

This function isn't working for me. Fix the bugs in this function

```
"""
def calculate_average(numbers):
    total = 0
    for i in range(len(numbers)):
        number = numbers[i]
        total += number
    average = total
rerun averave
```

```
def calculate_average(numbers):
    total = 0
    for i in range(len(numbers)):
        number = numbers[i]
        total += number
    average = total / len(numbers)
    return average
```

There were a few issues with the original code:

1. The variable "average" was assigned the value of "total" instead of the actual average.
2. There was a typo in the last line - "rerun" should be "return".
3. The average was not actually calculated - the corrected code divides the total by the length of the "numbers" list to get the average.

## Improve performance

- Identify where improvements can be made
- Models able to optimize and refactor

### Prompt:

Can this function be more efficient?

"""

```
def sum_of_n(n):
    result = 0
    for i in range(1, n+1):
        result += i
    return result
```

We can improve this function it by using a mathematical formula instead.

Here's the more efficient version:

```
def sum_of_n(n):
    return n * (n + 1) // 2
```

## Exercise: Generate and improve code with Azure OpenAI Service



Use the hosted lab environment if provided, or view the lab instructions at the link below:

<https://aka.ms/mslearn-openai-code-gen>

# Knowledge check



- 1 What are some benefits of using Azure OpenAI to generate code?**
  - Increase in efficiency and productivity
  - Increase in bugs and readability
  - Increase in time spent coding

---
- 2 What is the purpose of providing more context to the Azure OpenAI model when completing partial code?**
  - Providing more context makes the model less accurate
  - Providing more context doesn't affect the accuracy of the model
  - The more context provided to the model, the more accurate the response likely is

---
- 3 What is an example of how Azure OpenAI models can change coding language?**
  - Azure OpenAI models can only translate code from one language to another if the code is written in C#
  - If you have code in one language, but need it in another, Azure OpenAI can translate that for you in several languages
  - Azure OpenAI models can only translate code from Python to C# or Java

# Generate images with Azure OpenAI Service



# What is DALL-E?

Generate images with a description

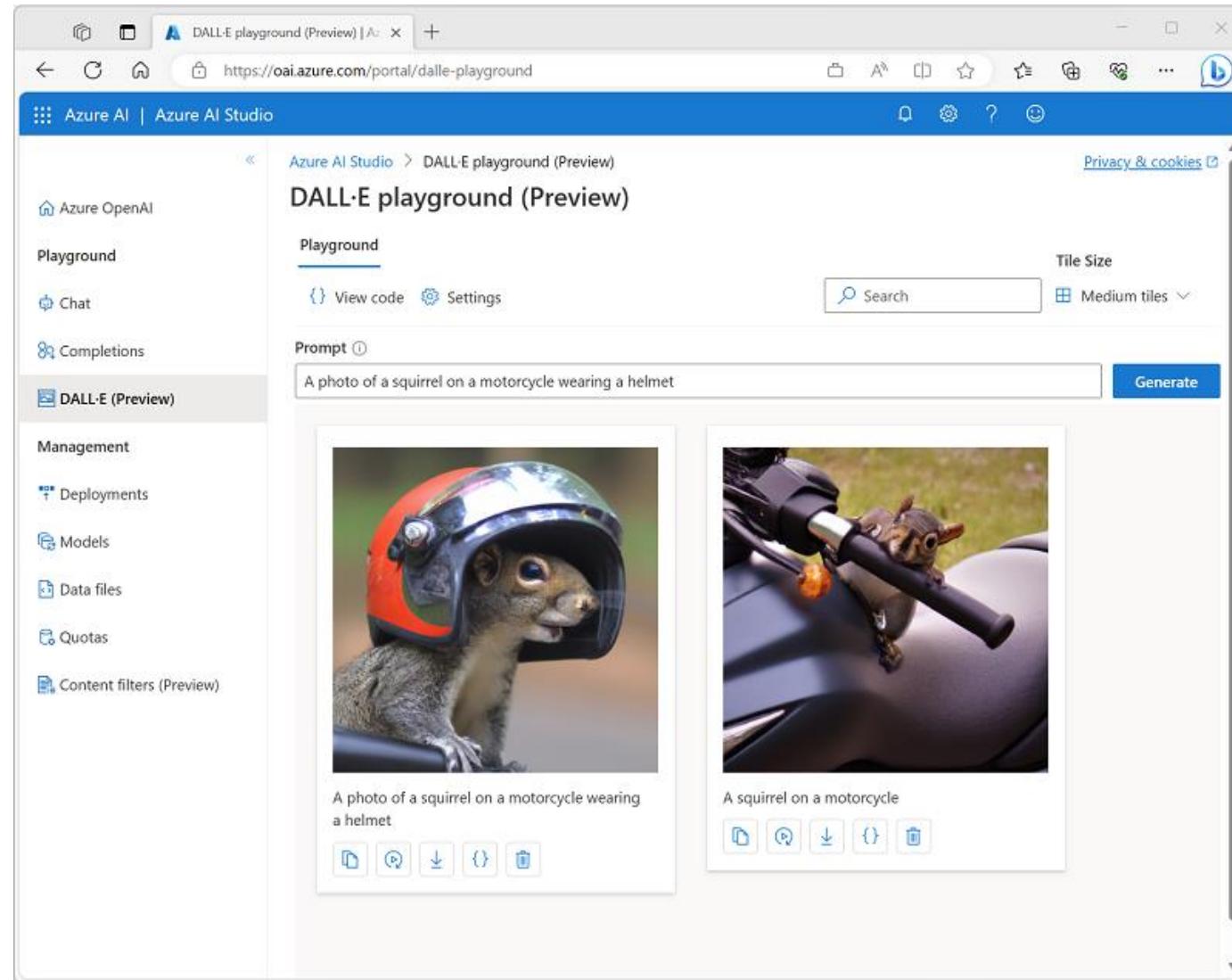
- Neural network based model for generating images
- Use natural language to describe what the image should be
- Specify content and style

**Prompt:**

a highland cow in a field on the coast of Scotland, digital art style



# Using DALL-E in Azure OpenAI Studio

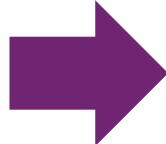


# Using the Azure OpenAI REST API

## DALL-E Endpoint

`https://endpoint.openai.azure.com/openai/images/generations:submit?api-version=api_version`

```
{  
    "prompt": "A badger wearing a tuxedo",  
    "n": 1,  
    "size": "512x512"  
}
```



```
{  
    "created": 1685130482,  
    "expires": 1685216887,  
    "id": "088e4742-89e8-4c38-9833-c294a47059a3",  
    "result":  
    {  
        "data":  
        [  
            {  
                "url": "<URL_to_generated_image>"  
            }  
        ]  
    },  
    "status": "succeeded"  
}
```

# Using the Azure OpenAI SDKs for DALL-E

Language specific SDKs are available for use in your applications, in both C# and Python.

Code structure follows a similar pattern for both languages.

The response content comes in the form of an image URL, for both REST and SDKs. The image can be downloaded or viewed directly from that URL.

Pseudo code structure:

```
<include library>

<create client>

<define image generation prompt and options>

<send request (can include options and prompt in request)>

<extract response content>
```

## Exercise - Generate images with a DALL-E model



Use the hosted lab environment if provided, or view the lab instructions at the link below:

<https://aka.ms/mslearn-openai-dall-e>

## Knowledge check



1 You want to use a model in Azure OpenAI to generate images. Which model should you use?

- DALL-E
  - GPT-35-Turbo
  - Text-Davinci
- 

2 Which playground in Azure OpenAI Studio should you use to explore image generation

- Completions
  - DALL-E
  - Chat
- 

3 In an API request to generate images, what does the n parameter indicate?

- The description of the desired image
- The number of images to be generated
- The size of the image to be generated