

OpenAI and ChatGPT for Enterprises – Introducing Azure OpenAI



In this chapter, we'll focus on the enterprise-level applications of OpenAI models and introduce the partnership between OpenAI and Microsoft and **Azure OpenAI (AOAI)** Service. We will go through the milestones and developments of Microsoft in the field of **artificial intelligence (AI)**, highlighting the journey that brought the Azure cloud into the game of OpenAI, and why this is a game-changer for large organizations. Finally, we will consider the topic of responsible AI and how to make sure your AI system complies with ethical standards.

In this chapter, we will discuss the following topics:

The history of the partnership between Microsoft and OpenAI and the introduction of AOAI Service

The role of the public cloud in the context of OpenAI models

Responsible AI

By the end of this chapter, you will have learned about the main features of AOAI Service and how it differs from the OpenAI models we've discussed so far. You will also be familiar with the partnership history between Microsoft and OpenAI, and why there was the need for OpenAI models to be deployed on an enterprise-scale infrastructure. Finally, you will understand Microsoft's continuous and long-lasting commitment toward responsible AI and how it is benefiting AOAI Service.

Technical requirements

The following are the technical requirements for this chapter:

An Azure subscription, which you can create for free here:
<https://azure.microsoft.com/free/cognitive-services>.

Access granted to Azure OpenAI in the desired Azure subscription. At the time of writing, access to this service is granted only by application. You can apply for access to Azure OpenAI by completing the form at
<https://aka.ms/oai/access>.

OpenAI and Microsoft for enterprise-level AI – introducing Azure OpenAI

Microsoft has a long history of investing in AI research and development, with a focus on building AI-powered tools and services that can be used by businesses and individuals to solve complex problems and improve productivity.

It also boasts a series of milestones in terms of achieving human parity in AI domains such as speech recognition (2017), machine translation (2018), conversational Q&A (2019), image captioning (2020), and natural language understanding (2021).

Definition

Human parity in AI refers to the point at which an AI system can perform a task or tasks at a level that is equal to or indistinguishable from a human being. This concept is often used to measure the performance of AI systems, especially in areas such as natural language understanding, speech recognition, and image recognition.

Achieving human parity in AI is considered a significant milestone as it demonstrates the AI's ability to effectively match human capabilities in a given domain.

In the next few sections, we are going to explore the research history and background of Microsoft in the domain of AI, to fully understand its journey toward their partnership with OpenAI and, finally, the development of AOAI Service.

Microsoft AI background

Early research in the field of AI traces back to the late 1990s when Microsoft established its **machine learning (ML)** and applied statistics groups. Starting with those, Microsoft started researching and experimenting with intelligent agents and

virtual assistants. In this case, the prototype was Clippy, a personal digital assistant for Microsoft Office:

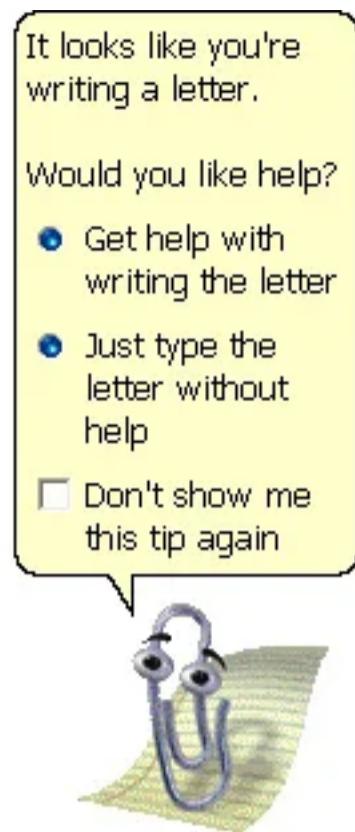


Figure 9.1 – Clippy, the default Office Assistant launched in 2000

Clippy was the forerunner of more sophisticated tools such as Cortana. Launched in 2014, Cortana is a digital assistant that uses **natural language processing (NLP)** and ML to provide personalized assistance to users.

Then, in 2016, as an expansion of Microsoft Project Oxford, Microsoft launched Microsoft Cognitive Services in the Azure cloud, a set of APIs that provide AI capabilities to developers without them requiring ML and data science expertise:

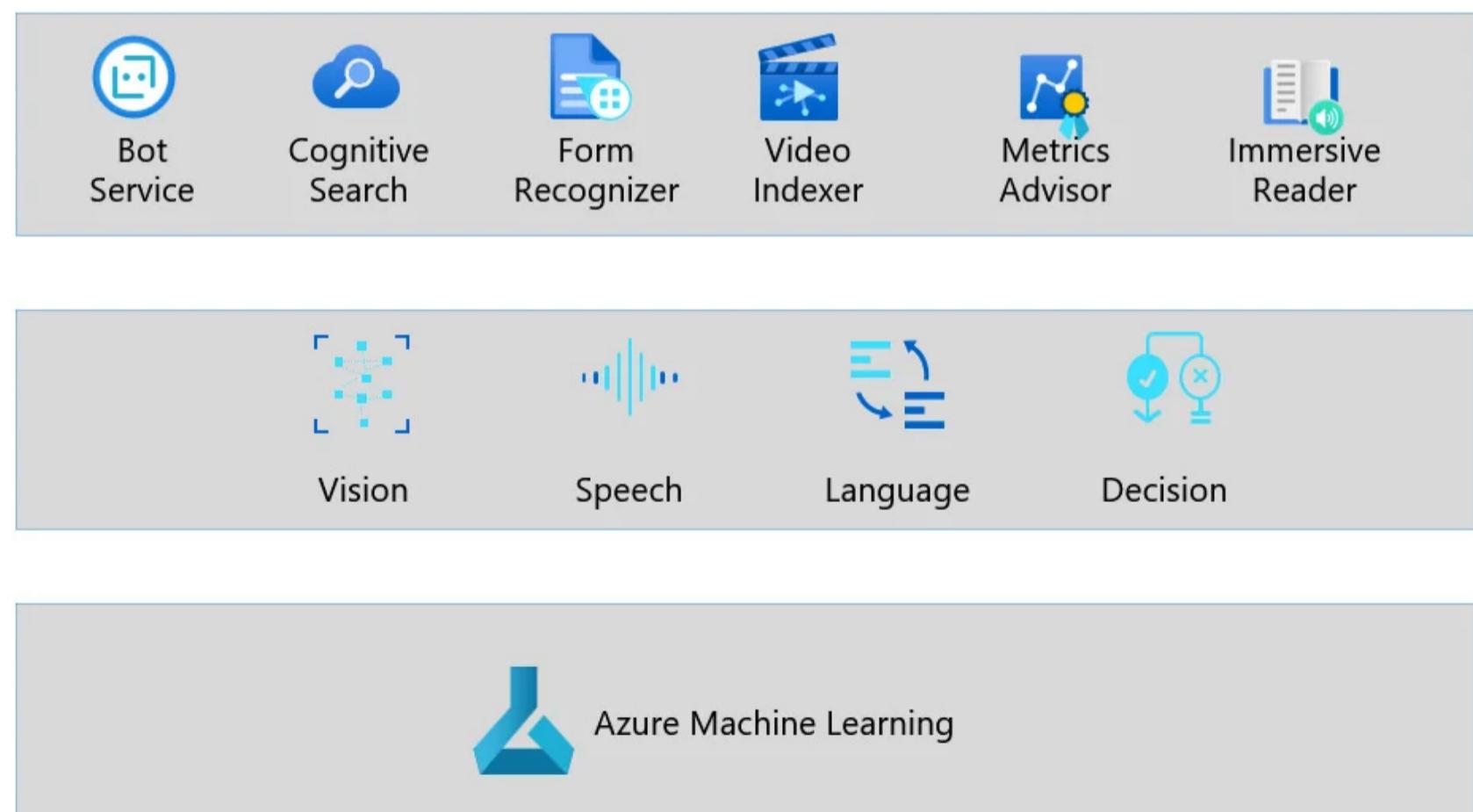


Figure 9.2 – Microsoft Azure AI services

With Cognitive Services, AI could finally be consumed by a wide range of users, from large enterprises to individual developers. From this, we witnessed what we now call **AI democratization**: AI is no longer a privilege for those who have deep technical knowledge and powerful and expensive hardware for model training. Cognitive Services has been developed for the following reasons:

So that anyone, from data scientists to business users, can leverage Cognitive Services with a no-code approach

To provide a set of pre-built models that have already been trained – that is, they are ready to use and don't need GPU-powered hardware to run

Microsoft's investments in AI can be seen from its acquisition of AI companies in recent years, including SwiftKey (a predictive keyboard app:

<https://blogs.microsoft.com/blog/2016/02/03/microsoft-acquires-swiftkey-in-support-of-re-inventing-productivity-ambition/>) in 2016, Maluuba (a deep learning startup: <https://blogs.microsoft.com/blog/2017/01/13/microsoft-acquires-deep-learning-startup-maluuba-ai-pioneer-yoshua-bengio-advisory-role/>) in 2017, and Bonsai (a platform for building AI models:

<https://blogs.microsoft.com/blog/2018/06/20/microsoft-to-acquire-bonsai-in-move-to-build-brains-for-autonomous-systems/>) in 2018.

Among the companies Microsoft invested in and partnered with, there is also OpenAI.

The partnership between the two tech companies began in 2016 when OpenAI agreed to leverage Microsoft's Azure cloud infrastructure to run its AI experiments. Later on, in 2019, Microsoft announced a \$1 billion partnership with OpenAI

(<https://news.microsoft.com/2019/07/22/openai-forms-exclusive-computing-partnership-with-microsoft-to-build-new-azure-ai-supercomputing-technologies/>) to develop AI models and technologies that can be used for the benefit of humanity.

This partnership is based on the following three main pillars:

Microsoft and OpenAI will jointly build new Azure supercomputing infrastructure to train AI models

OpenAI will make its models and technologies consumable from the Azure cloud

Microsoft will become OpenAI's preferred partner for commercializing new AI solutions to the market

Since then, the two companies kept investing and researching, and finally, in January 2023, a set of OpenAI models was made available in Azure via AOAI Service.

With the general availability of AOAI Service, a new milestone was reached and the Microsoft AI portfolio has been extended with the powerful large language models of OpenAI.

Azure OpenAI Service

AOAI Service is a product of Microsoft that provides REST API access to OpenAI's powerful language models such as GPT-3.5, Codex, and DALL-E. You can use these models for the very same tasks as OpenAI models, such as content generation, summarization, semantic search, natural language, and code translation.

In the context of the Microsoft Azure AI portfolio, AOAI Service is collocated among the following Cognitive Services offerings:

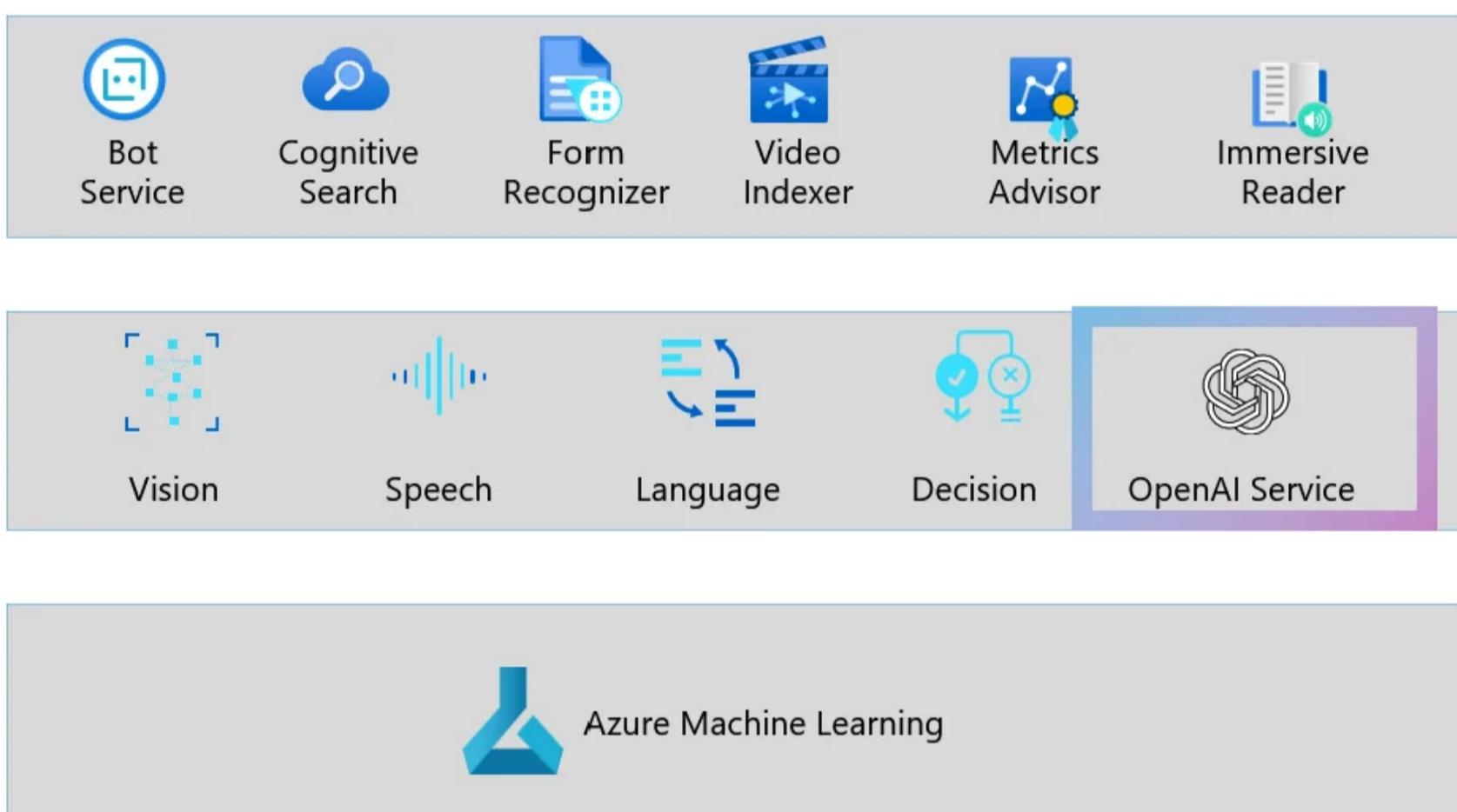


Figure 9.3 – AOAI Service General Availability (GA)

As with any other Cognitive Services offering, AOAI offers models that have already been trained and are ready to be consumed.

To create your AOAI resource, follow these instructions:

1. Navigate to the Azure portal at <https://ms.portal.azure.com>.
2. Click on **Create a resource**.
3. Type **azure openai** and click on **Create**.
4. Fill in the required information and click on **Review + create**.

This is shown in the following screenshot:

Azure services

Home > Create a resource >

Marketplace ...**Get Started**

Service Providers

Management

Private Marketplace

Private Offer Management

My Marketplace

Favorites

Recently created

Private products

Categories

AI + Machine Learning (1)

Analytics (0)

Home > Create a resource > Marketplace >

Create Azure OpenAI ...

azure openai

Azure benefit eligible only

Showing 1 to 1 of 1 results for 'azure openai'. [Clear search](#)

Create

Basics Tags Review + create

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.

Project details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription *

Valentina Alto

Resource group *

Create new

Instance details

Region *

South Central US

Name *

aoaitestvalentina

Pricing tier *

Review + create < Previous Next : Tags >

Figure 9.4 – Steps to create an AOAi resource

This process might take a few minutes. Once it is ready, you can directly jump to its user-friendly interface, AOAi Playground, to test your models before deploying them:

The image shows two side-by-side screenshots of the Azure OpenAI Studio interface.

Top Screenshot (Completions playground):

- Left Sidebar:** Includes links for Azure OpenAI, Playground, Chat, Completions (selected), Management, Deployments, Models, and Data Files.
- Header:** Shows "Azure OpenAI Studio > Completions playground" and "Privacy & cookies".
- Deployment Selection:** "Deployments" dropdown set to "text-davinci-003" and "Examples" dropdown set to "Generate an email". A "View code" button is also present.
- Text Input Area:** A large text input field with placeholder "Start typing here".
- Buttons at the bottom:** "Generate", "Undo", "Regenerate", and "Tokens: 0".
- Right Panel (Parameters):** Various AI generation parameters with sliders and input fields:
 - Temperature: 1
 - Max length (tokens): 350
 - Stop sequences
 - Top probabilities: 1
 - Frequency penalty: 0
 - Presence penalty: 0
 - Best of: 1
 - Pre-response text: Enter text
 - Post-response text: Enter text

Bottom Screenshot (Chat playground (Preview)):

- Left Sidebar:** Same as the top screenshot.
- Header:** Shows "Azure OpenAI Studio > Chat playground (Preview)" and "Privacy & cookies".
- Assistant setup (Left Panel):**
 - "Save changes" button.
 - "Specify how the chat should act": "Use a template to get started, or just start writing your own system message below. Want some tips? [Learn more](#)".
 - "Use a system message template": "Select a template" dropdown.
 - "System message": "You are an AI assistant that helps people find information."
 - "Examples": "Add examples to show the chat what responses you want. It will try to mimic any responses you add here so make sure they match the rules you laid out in the system message." with a "+ Add an example" button.
- Chat session (Middle Panel):**
 - "Clear chat", "View code", and "Show raw JSON" buttons.
 - Message history between "aoci" (bot) and a user:
 - Bot: "As an AI language model, I am not sure what you mean by "aoci." Could you please provide more context or clarify your question?"
 - User: "hello"
 - Bot: "We're experiencing heavy traffic right now. Please try again at a later time."
 - "User message" input field with placeholder "Type user query here. (Shift + Enter for new line)".
- Parameters (Right Panel):**
 - Deployment:** Set to "gpt-35-turbo".
 - Session settings:** "Past messages included: 10".
 - Current token count:** 100/4000.

Figure 9.5 – AOAI UI and Playground

Note that AOAI Playground looks almost identical to the OpenAI Playground version we saw in [Chapter 2](#). The difference here is that, to use AOAI models, you have to initiate a deployment, which is a serverless compute instance you can attach to a model. You can do so either in Playground or on the resource backend page in the Azure portal:

A

Deployments

Deployments enable you to make completions and search calls against a provided base model or your fine-tuned model. You can also scale up and down your deployments easily.

Create new deployment

Deployment name	Model name
text-davinci-003	text-davinci-003
code-search-babbage-text-001	code-search-babbage-text-001
gpt-35-turbo	gpt-35-turbo

B

aoaitest3 | Model deployments

Set up a model deployment to make API calls against a provided base model or a custom model. Finished model deployments are available for use. Your model deployment status will move to succeeded when the model deployment is complete and ready for use.

Create Model deployment

Model deployment name *

Model

Version

Save

Figure 9.6 – Creating a new AOAI deployment via Playground (A) or in the Azure portal (B)

For example, I created a deployment called **text-davinci-003** with an associated **text-davinci-003** model:

Deployments

Deployments enable you to make completions and search calls against a provided base model or your fine-tuned model. You can also scale up and down your deployments easily.

Create new deployment	Edit deployment	Delete deployment	Column options	Refresh	Open in Playground	Search	
Deployment name	Model name	M...	Sc...	Sc...	Sta...	Model dep...	Created at
text-davinci-003	text-davinci-003	1	Stand...	-	<input checked="" type="checkbox"/> ...	9/30/2024	3/16/2023 8:5...

Figure 9.7 – An active deployment of AOAI

In OpenAI Playground, we can test those models either directly via the user interface or by embedding their APIs into our applications. In the next section, we are going to explore how to interact with Playground and try different models' configurations. In [**Chapter 10**](#), we will learn how to integrate AOAI's Models API into enterprise applications.

Exploring Playground

AOAI Playground is the easiest way to get familiar with the underlying models and start planning which model's version is the most suitable for your projects. The user interface presents different tabs and workspaces, as shown in the following screenshot:

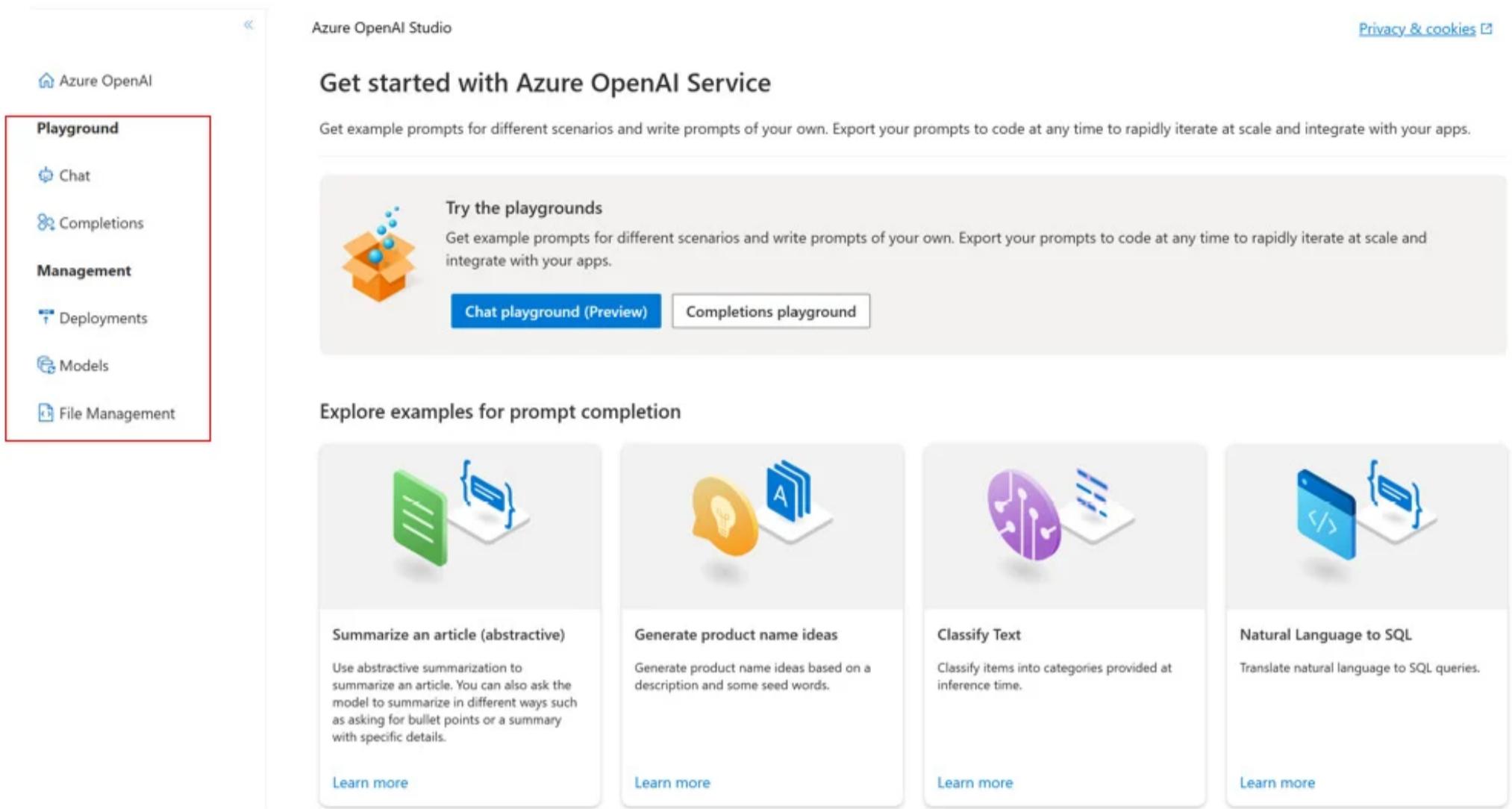


Figure 9.8 – Overview of AOAI Playground

Let's explore each of them:

Playground | Chat: The Chat workspace is designed to be only used with conversational models such as GPT-3.5-turbo (the model behind ChatGPT):

Playground

[Chat](#)[Completions](#)[Management](#)[Deployments](#)[Models](#)[File Management](#)

Assistant setup

Load example setup

Default

System message

You are an AI assistant that helps people find information.

> Few-shot examples

Chat session

Start chatting

Test your assistant by sending queries below. Then adjust your assistant setup to improve the assistant's responses.

User message

Type user query here. (Shift + Enter for new line)

Parameters

Deployments

gpt-35-turbo

Max response 800

Temperature 0.5

Top P 0.95

Stop sequence

Stop sequences

Frequency penalty 0

Presence penalty 0

[Learn more](#)

Figure 9.9 – AOA Chat workspace

It offers a similar experience to ChatGPT itself, with the possibility to configure your model with additional parameters (as we saw in [Chapter 2](#) with OpenAI Playground). Furthermore, there is an additional feature that makes the Chat workspace very interesting, known as **System message**:

Playground

[Chat](#)[Completions](#)[Management](#)[Deployments](#)[Models](#)[File Management](#)

Assistant setup

Load example setup

Default

System message

You are an AI assistant that helps people find information.

> Few-shot examples

Chat session

Start chatting

Test your assistant by sending queries below. Then adjust your assistant setup to improve the assistant's responses.

User message

Type user query here. (Shift + Enter for new line)

Parameters

Deployments

gpt-35-turbo

Max response 800

Temperature 0.5

Top P 0.95

Stop sequence

Stop sequences

Frequency penalty 0

Presence penalty 0

[Learn more](#)

Figure 9.10 – Example of System message

System message is the set of instructions we give the model to tell it how to behave and interact with us. As for the prompt, **System message** represents a key component of a model's configuration since it massively affects model performance.

For example, let's instruct our model to behave as a JSON formatter assistant:

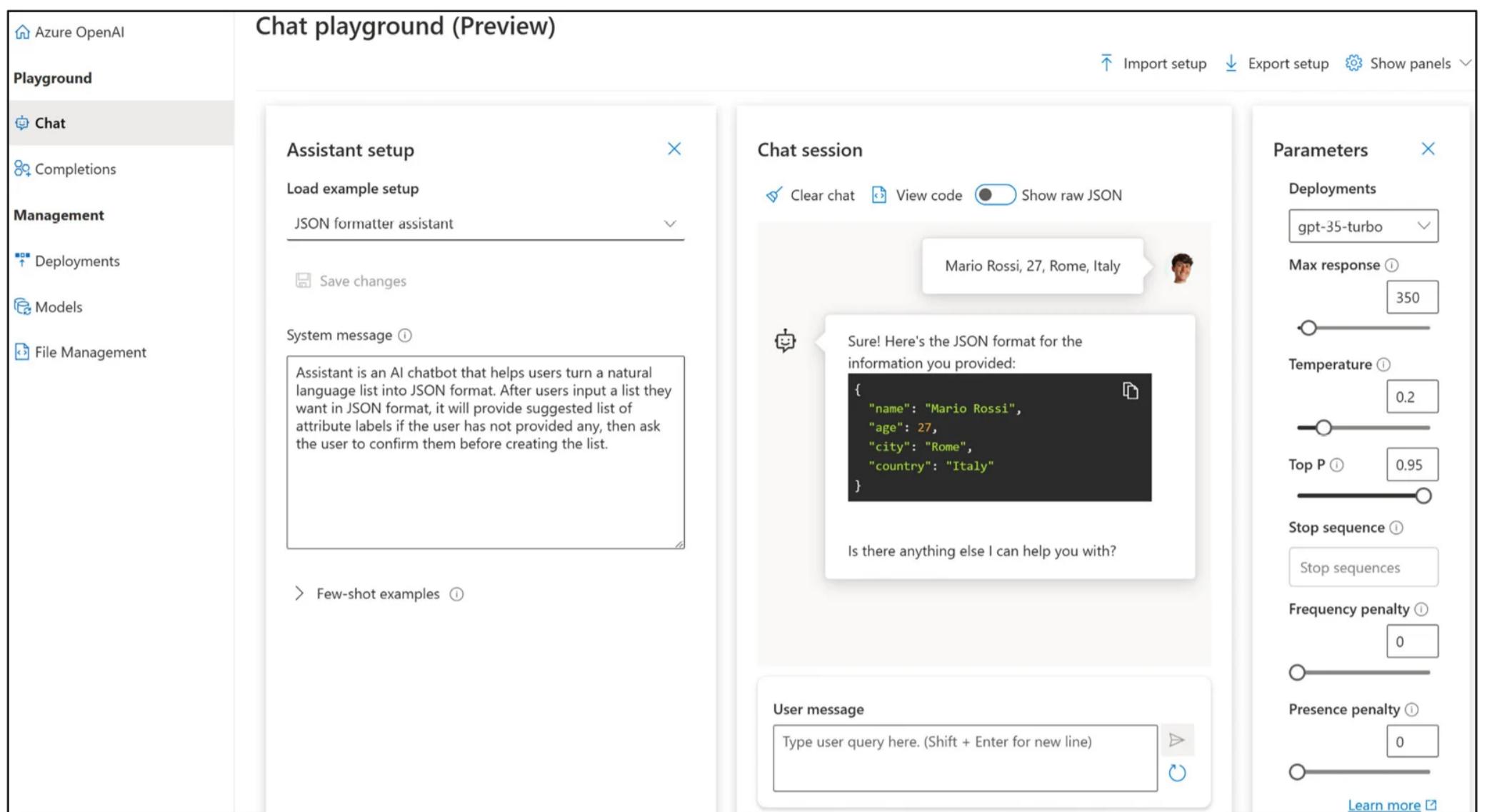


Figure 9.11 – Example of a model acting as a JSON formatter assistant

As you can see from the previous screenshot, the model was able to suggest a JSON file through some simple data, such as name and age, without the need to specify any labels.

Playground | Completions: Different from the previous workspace, the **Completions** workspace offers a sort of *white paper* where you can interact with your models. While GPT-3.5-turbo is designed for conversational tasks (which means it can be consumed via a chatbot-like interface), the GPT-3 series contains more general-purpose models and can be used for a wide range of language tasks, such as content generation, summarization, and so on.

For example, we could ask our model to generate a quiz by giving it a description of the topic and a one-shot example, as shown here:

The screenshot shows the Azure OpenAI Studio interface. On the left, a sidebar lists 'Azure OpenAI', 'Playground' (selected), 'Chat', 'Completions' (selected), 'Management', 'Deployments', 'Models', and 'File Management'. The main area is titled 'Completions playground' and has tabs for 'Deployments' (set to 'text-davinci-003') and 'Examples' (set to 'Generate a quiz'). A text input field says 'Generate a multiple choice quiz from the text below. Quiz should contain at least 5 questions. Each answer choice should be on a separate line, with a blank line separating each question.' Below this is a paragraph about neutron stars. To the right, a large section displays a generated quiz with five questions. The first question is 'Q1. What is a neutron star?' with options A, B, C, and D. The second question is 'Q2. What is the mass of a neutron star?' with options A, B, C, and D. The third question is 'Q3. What is the radius of a neutron star?' with options A, B, C, and D. The fourth question is 'Q4. What happens when a massive star explodes?' with options A, B, C, and D. The fifth question is 'Q5. What is the density of a neutron star?' with options A, B, C, and D. On the far right, a 'Parameters' panel includes sliders for 'Temperature' (0.8), 'Max length (tokens)' (500), and other settings like 'Stop sequences', 'Top probabilities', 'Frequency penalty', 'Presence penalty', 'Best of', 'Pre-response text', and 'Post-response text'. A 'Learn more' link is also present.

Figure 9.12 – Example of a GPT model generating a quiz

Finally, as per the **Chat workspace**, with **Completions**, you can configure parameters such as the maximum number of tokens or the temperature (refer to [Chapter 2](#) for a comprehensive list of those parameters and their meanings).

Management | Deployments: Within the **Deployments** tab, you can create and manage new deployments to be associated with AOAI models. They are depicted here:

The screenshot shows the Azure OpenAI Studio interface. The sidebar is identical to Figure 9.12. The main area is titled 'Deployments' and has a sub-header 'Deployments enable you to make completions and search calls against a provided base model or your fine-tuned model. You can also scale up and down your deployments easily the scale unit.' Below this is a table of deployments:

Deployment name	Model name	M...	Sc...	Sc...	Sta...	Model dep...	Created at
test1	text-davinci-002	1	Stand...	-	✓ ...	1/1/2024	11/15/2022 5:...
ada	text-ada-001	1	Stand...	-	✓ ...	3/1/2024	2/8/2023 10:5...
chatgpt	gpt-35-turbo	0301	Stand...	-	✓ ...	8/1/2023	3/9/2023 5:59 ...
code	code-davinci-002	1	Stand...	-	✓ ...	7/11/2024	2/1/2023 4:41 ...
embedding	text-embedding-ada-002	1	Stand...	-	✓ ...	2/2/2025	3/13/2023 6:1...
text-davinci-003	text-davinci-003	1	Stand...	-	✓ ...	9/30/2024	3/13/2023 5:1...

Figure 9.13 – List of AOAI deployments

Each deployment can host only one model. You can edit or delete your deployments at any time. As we mentioned previously, a model deployment is the enabler step for using either the **Completions** or **Chat** workspace within AOA Service.

Management | Models: Within this tab, you can quickly assess the models that are available within AOA Service and, among them, those that can be deployed (that is, a model that hasn't been deployed yet). For example, let's consider the following screenshot:

The screenshot shows the 'Models' section of the Azure OpenAI Studio. On the left, there is a sidebar with links: Azure OpenAI, Playground, Chat, Completions, Management, Deployments, Models (which is selected and highlighted in grey), and File Management. The main content area has a header 'Models' and a sub-header 'Provided models'. It includes a search bar and buttons for 'Deploy model', 'Create customized model', 'Column options', and 'Refresh'. A table lists nine models with columns for Model name, Model version, Created at, Status, and Deployable. The 'Deployable' column uses green checkmarks for Yes and red crossed-out checkmarks for No. The models listed are: text-similarity-curie-001, text-similarity-ada-001, text-embedding-ada-002, text-embedding-ada-002, text-davinci-003, text-davinci-002, text-curie-001, text-ada-001, and gpt-35-turbo.

Model name	Model version	Created at	Status	Deployable
text-similarity-curie-001	1	5/20/2022 2:00 AM	Succeeded	Yes
text-similarity-ada-001	1	5/20/2022 2:00 AM	Succeeded	Yes
text-embedding-ada-002	1	2/2/2023 1:00 AM	Succeeded	No ⓘ
text-embedding-ada-002	2	4/3/2023 2:00 AM	Succeeded	Yes
text-davinci-003	1	9/30/2022 2:00 AM	Succeeded	No ⓘ
text-davinci-002	1	1/22/2022 1:00 AM	Succeeded	No ⓘ
text-curie-001	1	3/1/2022 1:00 AM	Succeeded	Yes
text-ada-001	1	3/1/2022 1:00 AM	Succeeded	No ⓘ
gpt-35-turbo	0301	3/9/2023 1:00 AM	Succeeded	No ⓘ

Figure 9.14 – List of AOA models

Here, we have **text-similarity-curie-001**. It doesn't have an associated deployment, so it can be deployed (as the **Deployable** column shows). On the other hand, **text-similarity-ada-002** already has a deployment, so it is not available anymore.

Within this tab, you can also create a custom model by following a procedure called fine-tuning. We explored this in [Chapter 2](#):

The screenshot shows the 'Models' section of the Azure OpenAI Studio with the 'Create customized model' dialog open. The dialog has a title 'Create customized model' and a sidebar with options: 'Base model' (selected), 'Training data', 'Validation data', 'Advanced options', and 'Review and train'. It also includes a 'Base model type' dropdown and a 'Model suffix' input field. The main content area shows the same list of models as Figure 9.14, with the 'Create customized model' button highlighted with a red box. The table data is identical to the one in Figure 9.14.

Model name	Model version	Created at	Status	Deployable
text-similarity-curie-001	1	5/20/2022 2:00 AM	Succeeded	Yes
text-similarity-ada-001	1	5/20/2022 2:00 AM	Succeeded	Yes

Figure 9.15 – Example of model fine-tuning

Starting from this guided widget, you can upload your training and validation data to produce a customized model, starting from a base model (namely, **text-davinci-002**), which will be hosted on a dedicated deployment.

Note

In [Chapter 2](#), we saw that the training dataset should align with a specific format of the following type (called JSONL):

```
{"prompt": "<prompt text>", "completion": "<ideal generated text>"}
```

```
{"prompt": "<prompt text>", "completion": "<ideal generated text>"}
```

```
{"prompt": "<prompt text>", "completion": "<ideal generated text>"}
```

...

To facilitate this formatting, OpenAI has developed a tool that can format your data into this specific format ready for fine-tuning. It can also provide suggestions on how to modify data so that the tool can be used for fine-tuning. Plus, it accepts various data formats as inputs, including CSV, TXT, and JSON.

To use this tool, you can initialize the **OpenAI command-line interface (CLI)** by running the following command:

```
pip install --upgrade openai
```

Once initialized, you can run the tool, as follows:

```
openai tools fine_tunes.prepare_data -f <LOCAL_FILE>
```

Management | File Management: Finally, within the **File Management** tab, you can govern and upload your training and test data directly from the user interface, as shown here:

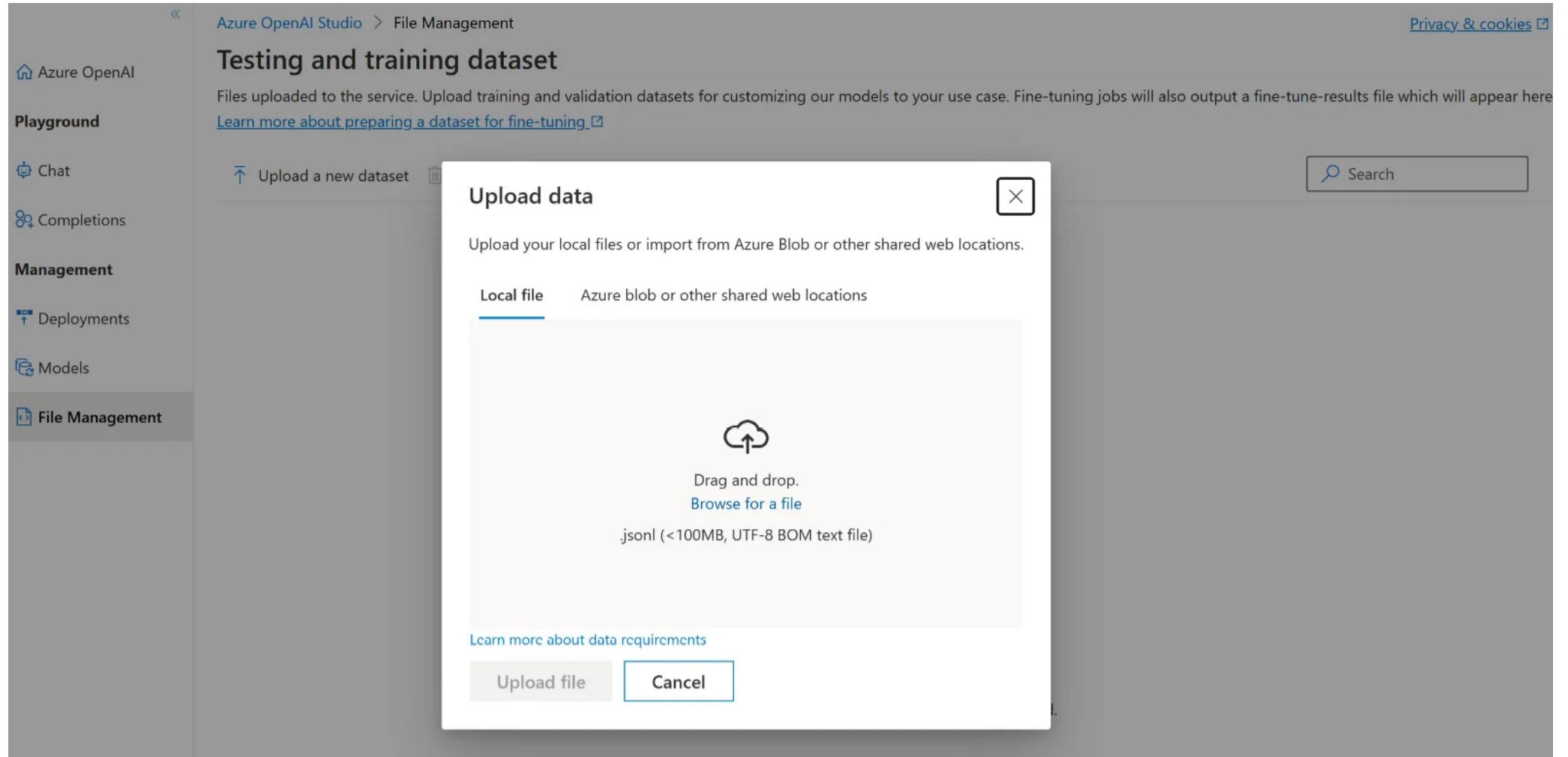


Figure 9.16 – Example of uploading a file within AOAI Service

You can decide to upload files by selecting **Local file** or **Azure blob or other shared web locations**.

Once you've uploaded your files, you will be able to select them while creating customized models, via the **Models** tab.

Finally, as mentioned in the previous section, each model comes with a REST API that can be consumed in your applications.

In the next chapter, we will see many end-to-end implementations of using AOAI's Models API. However, before we jump into that, we need to understand how AOAI differs from the standard OpenAI models and why the Azure cloud became part of the game.

Why introduce a public cloud?

At the beginning of this chapter, we saw how Microsoft and OpenAI have partnered in recent years and how Microsoft's cloud, Azure, became the *gym* for OpenAI model training. However, it also became the cloud infrastructure where OpenAI models can be consumed.

But what is the difference between using models from OpenAI and Azure OpenAI? The difference is the underlying infrastructure: with Azure OpenAI, you are leveraging your *own* infrastructure while living in your *own secured* subscription. This brings a series of advantages:

Scalability and flexibility: You can benefit from the scalability of Azure and accommodate the elastic usage of AOA models. From small pilots to enterprise-level production projects, AOA allows you to leverage the required capacity and scale up or down if necessary.

Security and compliance: You can use role-based authentication and private network connectivity to make your deployment more secure and trusted. You can also train your AI model while having full control of your data.

Regional availability: You can run your AI workloads on the Azure global infrastructure that meets your production needs.

Built-in responsible AI: You can use content filtering to ensure that your AI model generates appropriate and ethical output.

With the OpenAI models available in Azure, we can elevate the game to the enterprise and production levels, meeting all security and capacity requirements typical of large organizations.

One of the previously mentioned benefits deserves a particular focus: responsible AI. The rapid development of AI technologies also needs to be addressed in terms of ethical tools. This is what Microsoft has been studying since 2016, as we will explore in the next section.

Understanding responsible AI

We mentioned the built-in responsible AI as one of the key features of AOA. However, to fully understand it, we first need to understand Microsoft's commitment and journey toward responsible AI.

Microsoft's journey toward responsible AI

Microsoft soon recognized that as AI technologies continue to advance and become more integrated into our lives, there is a growing need to ensure that those systems are developed and used responsibly, ethically, and in ways that benefit everyone.

The beginning of this journey traces back to 2016 when Microsoft's CEO Satya Nadella penned an article exploring how humans and AI can work together to solve society's greatest challenges and introducing the first concepts of responsible AI, among which are transparency, fairness, and that it is designed for privacy and to assist humanity.

Shortly after, in 2017, Microsoft formalized those concepts with the first AI ethics committee – **Aether** (short for **AI, Ethics, and Effects in Engineering and Research**) – formed as an advisory group for the Microsoft senior leadership team.

AETHER spent time listening to customers and internal experts, and then partnered with legal affairs to publish the book titled *The Future Computed: Artificial Intelligence and its role in society* in January 2018. In this book, Microsoft identified six principles meant to guide a company's development of AI systems, as well as to help inform the broader industry and society as a whole about responsible AI practices.

Microsoft's six principles for responsible AI are as follows:

Fairness: Microsoft aims to create AI systems that are unbiased and treat all individuals and groups fairly, without discrimination or prejudice

Reliability and safety: Microsoft seeks to create AI systems that are robust, reliable, and secure, and that do not compromise safety or create unintended harm

Privacy and security: Microsoft values the privacy and security of individuals and their data, and works to protect them through transparency and responsible use of AI technologies

Inclusiveness: Microsoft believes that AI should be designed to empower and include individuals from diverse backgrounds, and to foster equal opportunities for all

Transparency: Microsoft believes in transparency and accountability for the decisions and actions of AI systems and is committed to providing clear explanations for their outcomes

Accountability: Microsoft accepts responsibility for the impact of its AI systems on society and the environment, and seeks to promote ethical and responsible practices in the development and use of AI

Microsoft follows these principles with the help of committees that offer guidance to its leadership, engineering teams, and every other team within the company.

Microsoft also has a **Responsible AI Standard** that provides a framework for building AI systems responsibly.

Following the publication of that book, Microsoft kept investing and researching in the following fields of responsible AI:

From the contribution to government regulation in the field of facial recognition (2018, <https://www.geekwire.com/2018/microsoft-calls-government-regulation-facial-recognition-technology/>, <https://blogs.microsoft.com/on-the-issues/2018/12/17/six-principles-to-guide-microsofts-facial-recognition-work/>) to the establishment of responsible AI in systems engineering or RAISE (2020, <https://www.microsoft.com/en-us/ai/our-approach?activetab=pivot1%3aprimaryr5>)

The development of responsible AI tools in the areas of ML interpretability, unfairness assessment and mitigation, error analysis, causal inference, and counterfactual analysis (2021, <https://responsibleaitoolbox.ai/>)

The following diagram shows the entire journey for responsible AI:

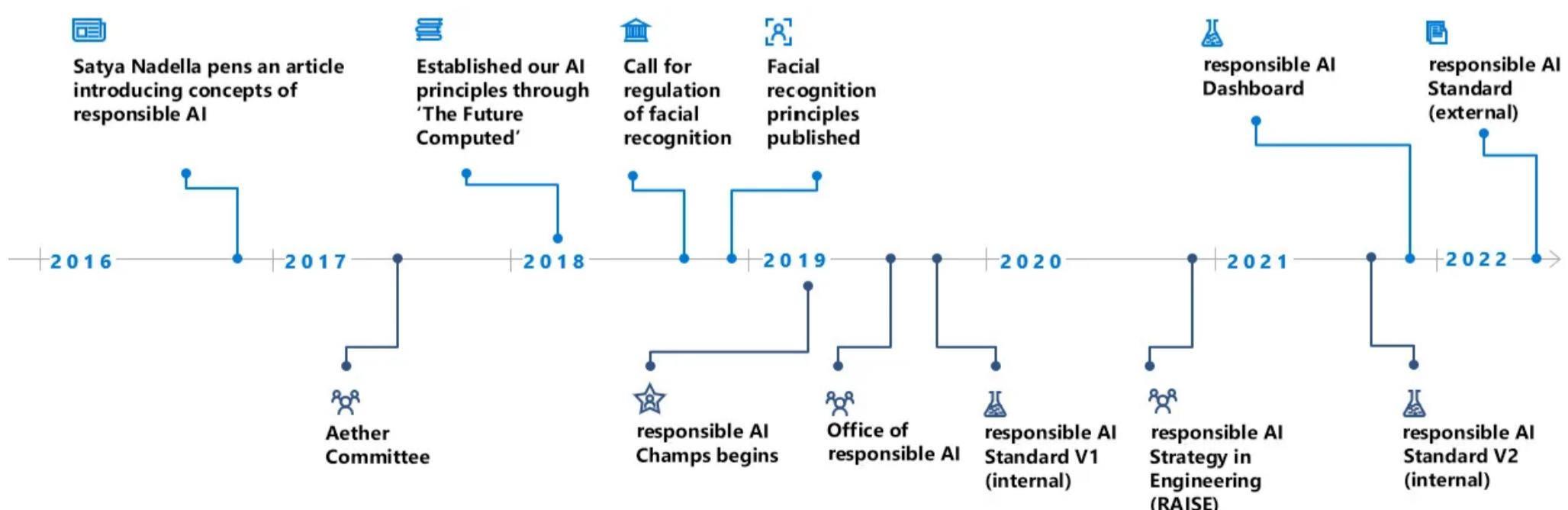


Figure 9.17 – Microsoft’s responsible AI journey

Microsoft’s commitment to responsible AI is reflected in the way its products are designed and the best practices and guidelines provided.

Of course, this also applies to AOAI Service. As we will see in the next section, AOAI Service comes with a built-in responsible AI at a different level.

Azure OpenAI and responsible AI

When it comes to AOAI Service, we can talk about responsible AI at the following two levels:

Built-in: In this case, we refer to AOAI embedded features of responsible AI that are enforced by a content management system. This system utilizes a series of classification models to detect harmful content. The system works alongside core models to filter content by analyzing both the input prompt and generated content. In cases where harmful content is identified, you'll receive either an error on the API call if the prompt was detected as inappropriate, or see that the `finish_reason` parameter on the response in JSON will be `content_filter` to signify that some of the generation was filtered.

Code of conduct and best practices: As for its other AI services, Microsoft provides **Transparency Notes** for AOAI. This application aims to promote an understanding of how AI technology works, its limitations and capabilities, and the importance of considering the entire system, including people and the environment. These notes can be used by developers and system owners to create AI systems that are fit for their intended purpose and, in the specific case of AOAI, help identify those scenarios that might trigger the built-in content filter.

Both the built-in capabilities and Transparency Notes are manifestations of Microsoft's effort to apply ethical AI practices in real-world scenarios, guided by their AI principles.

In conclusion, as responsible AI for Microsoft signifies the company's unwavering commitment to ethical AI development and deployment, AOAI also benefits from this commitment.

Summary

In this chapter, we saw how the partnership between OpenAI and Microsoft has brought about a powerful and innovative AI solution for enterprise-level organizations: AOAI. This service combines OpenAI's cutting-edge technology with

Microsoft's extensive cloud infrastructure to provide businesses with a scalable and customizable platform for building and deploying advanced AI applications.

We also dwelled on Microsoft's strong focus on responsible AI practices and ethics, and how AOA Service reflects this commitment to responsible AI, with features such as a content filter built into the platform.

As AI continues to transform industries and shape our future, the collaboration between OpenAI and Microsoft marks an important milestone in the development of enterprise-level AI solutions. AOA empowers businesses to harness the power of AI to drive growth and innovation while ensuring ethical and responsible practices.

In the next chapter, we will dive deeper into concrete use cases that enterprises are developing with the AOA Models API. We will also see an end-to-end implementation of a potential use case that uses Python and Streamlit so that you can experience firsthand how AOA's models can infuse your applications with AI.

References

<https://blogs.microsoft.com/blog/2023/01/23/microsoftandopenaiextendpartnership/>

<https://news.microsoft.com/2019/07/22/openai-forms-exclusive-computing-partnership-with-microsoft-to-build-new-azure-ai-supercomputing-technologies/>

<https://azure.microsoft.com/en-us/blog/general-availability-of-azure-openai-service-expands-access-to-large-advanced-ai-models-with-added-enterprise-benefits/>

<https://slate.com/technology/2016/06/microsoft-ceo-satya-nadella-humans-and-a-i-can-work-together-to-solve-societys-challenges.html>

<https://www.geekwire.com/2018/microsoft-calls-government-regulation-facial-recognition-technology/>

<https://blogs.microsoft.com/on-the-issues/2018/12/17/six-principles-to-guide-microsofts-facial-recognition-work/>

<https://www.microsoft.com/en-us/ai/our-approach?activetab=pivot1%3aprimaryr5>

<https://responsibleaitoolbox.ai/>

<https://www.microsoft.com/en-us/research/publication/human-parity-on-commonsenseqa-augmenting-self-attention-with-external-attention/>

<https://learn.microsoft.com/en-gb/azure/cognitive-services/openai/how-to/fine-tuning?pivots=programming-language-studio#openai-cli-data-preparation-tool>

[Previous Chapter](#)

[Next Chapter](#)