

# **SISTEMAS COMPUTACIONAIS AVANÇADOS (SISTCA) ADVANCED COMPUTING SYSTEMS**

Degree in Telecommunications and Informatics Engineering

Instituto Superior de Engenharia do Porto, Politécnico do Porto

**Lab Classes Script:**  
**OpenAI**

Version	Date	Authors	Update information
V1.0	April 2024	Patrícia Sousa (1210713), Carlos Alves (1211604), José Leal(1211066), Tiago Ribeiro (1210924) Supervisor: Paula Viana (pmv)	Original version

Table 1: Version Control

# Contents

<b>1</b>	<b>Introduction</b>	<b>6</b>
1.1	Context . . . . .	6
1.2	Motivation . . . . .	6
1.3	Objectives . . . . .	6
1.4	Document Structure . . . . .	7
<b>2</b>	<b>Theoretical (scientific/technological background)</b>	<b>8</b>
2.1	State-of-the-art . . . . .	8
2.2	OpenAI Features . . . . .	8
2.2.1	GPT . . . . .	8
2.2.2	Other models . . . . .	9
<b>3</b>	<b>Setup/Installation</b>	<b>10</b>
3.1	Creating an OpenAI Account . . . . .	10
3.2	Setting up Your Development Environment . . . . .	11
3.2.1	Windows . . . . .	11
3.2.2	Linux . . . . .	12
3.2.3	MacOS . . . . .	14
3.3	Installing Jupyter . . . . .	15
3.4	Installing required packages . . . . .	15
<b>4</b>	<b>Tutorial/Functionality</b>	<b>16</b>
4.1	Chat Completions . . . . .	16
4.2	Assistants . . . . .	17
4.2.1	Assistant . . . . .	17
4.2.2	Function Calling . . . . .	20
4.2.3	Code Interpreter . . . . .	27
4.2.4	File Search . . . . .	30
4.3	Embeddings . . . . .	32
4.4	Vision . . . . .	34
4.5	Image Generation (DALL-E) . . . . .	34
4.5.1	Image Generation . . . . .	35
4.5.2	Edit (DALLE 2 only) . . . . .	36
4.5.3	Variations (DALL-E 2 only) . . . . .	37
4.6	TTS . . . . .	37
4.7	Whisper . . . . .	38
4.8	Moderation . . . . .	40
<b>5</b>	<b>Exercises</b>	<b>42</b>
5.1	Exercise A . . . . .	42
5.2	Exercise B . . . . .	47
<b>6</b>	<b>Challenge</b>	<b>49</b>
<b>7</b>	<b>Future Work</b>	<b>50</b>
<b>8</b>	<b>Appendix</b>	<b>51</b>
8.1	Exercise A - Solution . . . . .	51
8.2	Exercise B - Solution . . . . .	55

## List of Tables

1	Version Control . . . . .	2
2	AI Companies. . . . .	8
3	TTS Information Table . . . . .	38

## List of Figures

1	OpenAI Website Access . . . . .	10
2	Select the API Option . . . . .	10
3	OpenAI Account Creation Confirmation . . . . .	11
4	Download Python . . . . .	11
5	Setup Python . . . . .	14
6	Run status flow diagram . . . . .	18
7	Resulting graph . . . . .	30
8	Image generated . . . . .	36
9	Input/Output from the edit endpoint . . . . .	36
10	Input/Output from the variation endpoint . . . . .	37

# 1 Introduction

## 1.1 Context

In this script we will present an API developed by openAI, one of the leading organisations in artificial intelligence (AI) research and development. Launched in 2020, it represents a significant milestone in accessing advanced AI technology. It offers an accessible and simplified interface for developers to integrate AI capabilities into a variety of applications, from text analysis to content generation.

The OpenAI API is based on recent advances in machine learning, in particular, deep neural network architectures, such as GPT (Generative Pre-trained Transformer) language models. These models have revolutionised the way machines understand and generate natural language, enabling tasks such as automatic translation, text summarising, question answering, and content generation.

GPT is trained with a vast amount of textual data from the internet, allowing it to capture the nuances and complexities of human language. Through the pre-training process, the model learns to represent knowledge in a general way, without the need for specific training for a particular task.

The OpenAI API leverages these pre-trained language models, providing a simple and effective interface for developers to take advantage of their capabilities. This means users can easily integrate advanced AI capabilities into their applications, without needing to understand complex implementation details or model training.

Additionally, it is designed with a focus on security and ethics, including measures to prevent malicious or harmful use of AI technology. This reflects OpenAI's commitment to promoting the responsible development of artificial intelligence and ensuring that its benefits are extremely accessible and used for the good of society.

Lastly, for a more interactive experience for the user, this script is available in another format, jupyter notebook. A free open-source software launched in 2014 by IPython that provides a notebook, this is, a shareable document that combines computer code, plain language descriptions, data, rich visualizations like 3D models, charts, graphs and figures, and interactive controls in one file, creating a fast and easy way to prototype and explain code, visualize data, and share ideas [1].

## 1.2 Motivation

One of the main motivations for choosing this topic lies in the importance of artificial intelligence in society today, some of whose main impacts include the automation and efficiency of tasks, more accurate medical diagnoses, learning adapted to the needs of each student, as well as the existence of virtual assistants.

By integrating artificial intelligence into various aspects of society, it is possible to solve complex problems, improve quality of life and drive innovation. However, it is essential that the development and implementation of AI is carried out in an ethical and responsible manner, taking into account privacy and security issues to ensure that the benefits are widely distributed and the risks mitigated.

In choosing this subject, it is recognised that there is a need to become familiar not only with the fundamental concepts of artificial intelligence, but also with the practical tools for its implementation. To this end, it is necessary not only to master theoretical principles, but also to explore programming languages and accessible platforms, such as the OpenAI API.

Whether in the SISTCA course or in other contexts, the OpenAI API offers a versatile and robust environment for exploring the potential of artificial intelligence.

## 1.3 Objectives

The main objective of this script is to explore the various features offered by the OpenAI API and learn how to integrate them, promoting a comprehensive and practical understanding of the platform's capabilities, providing users with a solid foundation to explore and use the following available tools (2.2):

- **Chat Completions**, to understand the generic implementation of chat GPT based assistant;
- **Assistants API + Tools** To create a program that provides the user with media suggestions such as books, films or TV shows;

- **Embeddings** Learn how to represent words, sentences, paragraphs, or entire documents in a continuous vector space;
- **Vision** The use of Vision functionality to extract information from an image.
- **Image Generation/DALL-E** Learn how to generate AI created images.
- **TTS (Text to Speech)** Turn strings into a robotised speech using AI;
- **Whisper** Capture an audio and translate it to different languages;
- **Moderation** Explore OpenAI's Moderation functionality to detect inappropriate content.

In addition, two practical exercises will be developed to apply the acquired knowledge and a challenge will be proposed to test the reader's understanding and ability to use the OpenAI API.

## 1.4 Document Structure

This tutorial is organised into seven chapters with the following structure: Introduction, Theoretical, Setup/Installation, Tutorial/Functionality, Exercises followed by a Challenge and ends with a closing chapter.

In the Introduction, some concepts about the OpenAI API will be presented that will be deepened throughout the script.

The theoretical part, will introduce the State-of-the-art and detail the API's Features.

Then, the third chapter focuses on Setup/Installation, explaining how to create an open AI account and set up a development environment.

In the Tutorial/Functionality part, a detailed tutorial of the available features will be provided.

After all these topics, two exercises will be developed and their corresponding resolution will be presented and a final challenge to test the knowledge acquired by the users.

After the tutorials two exercises are proposed, which are followed by a final challenge.

Finally, we conclude by discussing future improvements to the script.

## 2 Theoretical (scientific/technological background)

Before diving into specific AI models it's important to have a general understanding of how AI generally works.

Most AI solutions in natural language processing are based on broader **Large language models (LLMs)** these are algorithms that have been trained on vast amounts of data with the purpose of understanding and generating human like text. LLMs make use of the **transformer** type of architecture, a deep learning technique that, in short, represents text via numerical representations known as tokens and gives them different weights so as to be able to contextualise words and find similarities [2].

One example of a transformer based model is the **Generative pre-trained transformer (GPT)**, an AI model that has been pre-trained on large sets of data via the use of the transformer architecture for general purpose tasks. Furthermore, these models can be fine tuned to achieve greater performance in more specific tasks, such is the case of image generation models [3].

Lastly, in order to make use of these models it's essential to understand the concept of prompts. That is the name given to the textual inputs given to the model. These inputs are then broken down into the aforementioned tokens via a process we call tokenization, which facilitates the use of the models language structure.

### 2.1 State-of-the-art

AI, Artificial Intelligence, refers to a simulation of human intelligence in machines programmed to mimic human cognitive processes and actions. The concept of AI is not new, it has been around since the mid-20th century, but in the last few years there have been significant advances for AI. With these advances, due to its potential, AI has become increasingly important across various industries such as healthcare, finance, manufacturing, education, amongst others. Nowadays, one of the most famous companies developing AI products is OpenAI, which we are going to base our work on. However, just like any other business, OpenAI has its competitors. Some of them currently only dispose of a chat bot, and even that, at the time, are not available in Portugal. Listed below are some companies in the AI business.

AI	OpenAI[4]	X[5]	Anthropic[6]	Deepmind[7]	Cohere[8]
ChatBot	ChatGPT	Grok	Claude	Gemini	Coral
API	Yes	No	Yes	Yes	Yes
Inputs	Text, Audio, Image, Video		Text	Text, Image	Text
Outputs	Text, Audio, Image, Video		Text	Text	Text
ChatBot and API Availability in Portugal	Both	None	API Only	ChatBot Only	Both

Table 2: AI Companies.

### 2.2 OpenAI Features

OpenAI's API offers a vast array of cutting-edge AI models based on deep learning and natural language processing techniques. These models have been trained on vast datasets and fine tuned to fulfil a variety of tasks such as text and image generation, audio and text conversions, amongst other things.

#### 2.2.1 GPT

The GPT (Generative Pre-trained Transformer) series is OpenAI's main set of large language models. These models are trained to understand and generate natural language text based on contextual inputs so as to better communicate with humans. The most widespread version, GPT-3.5 currently serves as the model that powers the free version of ChatGPT. Its understanding of human language allows for coherent conversations which makes it a suitable chat bot. GPT-4 improves upon its predecessor with a smarter and more knowledgeable model that provides greater accuracy across various tasks. In particular, GPT-4 introduces Vision as a new feature, that enables it to process image inputs, making it useful in a wider range of applications.

### **Function calling**

The GPT based models are capable of calling previously specified functions in response to user actions or prompts by calling external APIs to retrieve data or to automate procedures like sending an email or extracting and sorting data from a document.

### **Assistants**

The AI Assistants functionality leverages the use of the GPT models alongside function calling and other tools like file retrieval and code interpreter to allow users to create custom assistants that fulfil more specific tasks based on the provided instructions.

#### **2.2.2 Other models**

##### **DALL-E**

The DALL-E model is capable of generating images from natural-language text descriptions, as well as modifying existing ones by feeding the model instructions from a text prompt.

##### **TTS**

The Text-To-Speech model is capable of converting text into a natural sounding speech. Note that it currently only supports the English language.

##### **Whisper**

Whisper does the opposite of the TTS model: it takes an audio input and then transcribes it into text. Unlike text to speech, Whisper is capable of understanding multiple languages, as such it can be used to identify the input language and translate the contents of the speech into English.

##### **Embeddings**

Text embeddings are vectorial representations of strings of text, such as words or phrases. By comparing two or more vectors we can infer their similarity. This mechanism is highly useful in applications such as search engines or product recommendations due to its ability of evaluating similarity between text strings.

##### **Moderation**

OpenAI's Moderation model is designed to verify if a certain piece of text includes any content that could be classified as hateful, violent, sexual, harmful or otherwise inappropriate. Whilst OpenAI's own use of the model aims to ensure that content complies with their usage policies [4], this model is suitable for any application that aims to ensure a safe digital environment.

### 3 Setup/Installation

#### 3.1 Creating an OpenAI Account

In this section, we will guide you through the process of setting up an OpenAI account. Whether you're a developer, researcher, or simply curious about AI, having your own account opens the door to the vast possibilities of artificial intelligence.

Firstly, navigate to the OpenAI website [4] to create or log into your account.

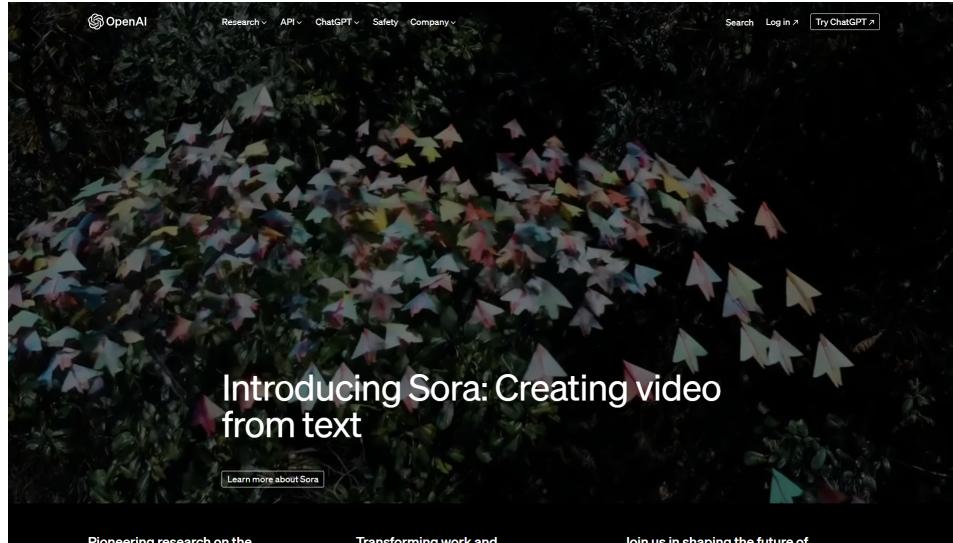


Figure 1: OpenAI Website Access

Upon logging in with your email, you will be presented with two options: ChatGPT and API. Select the API option to access the documentation.

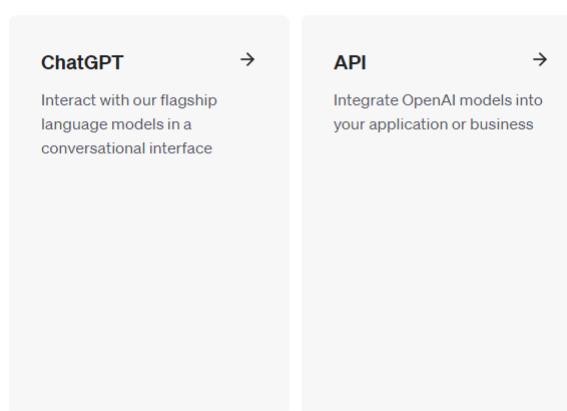


Figure 2: Select the API Option

Congratulations! You have now successfully created an operational OpenAI account.

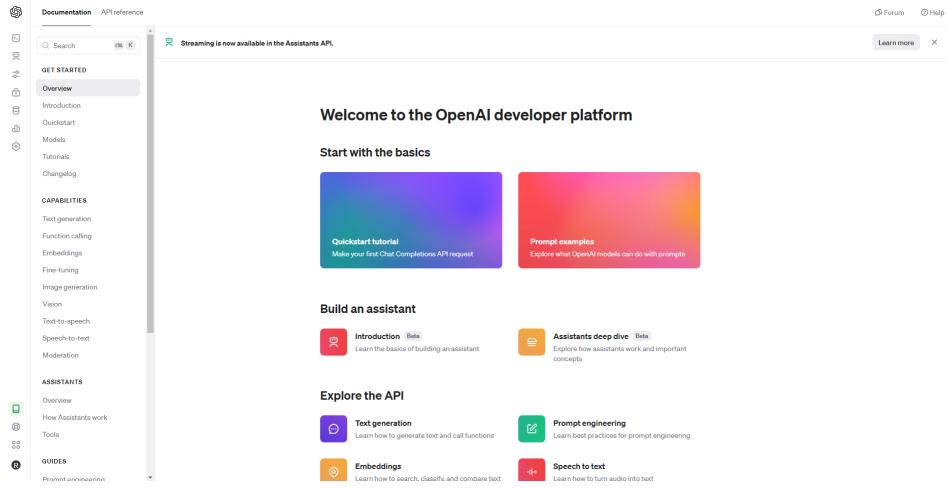


Figure 3: OpenAI Account Creation Confirmation

## 3.2 Setting up Your Development Environment

Setting up a proper development environment is crucial for working efficiently with AI applications. Ensure that you have the necessary tools and libraries installed on your system. For most AI development tasks, Python is the recommended programming language due to its extensive ecosystem of AI and machine learning libraries. So, that is exactly what we are going to help you with in this sub-section.

### 3.2.1 Windows

The first step is accessing python's official website and downloading it. In case you are not sure if you have already installed Python in the past, just type "cmd" in your search-bar and then type "python". If you are having trouble installing, maybe try checking Python's beginners guide [9].

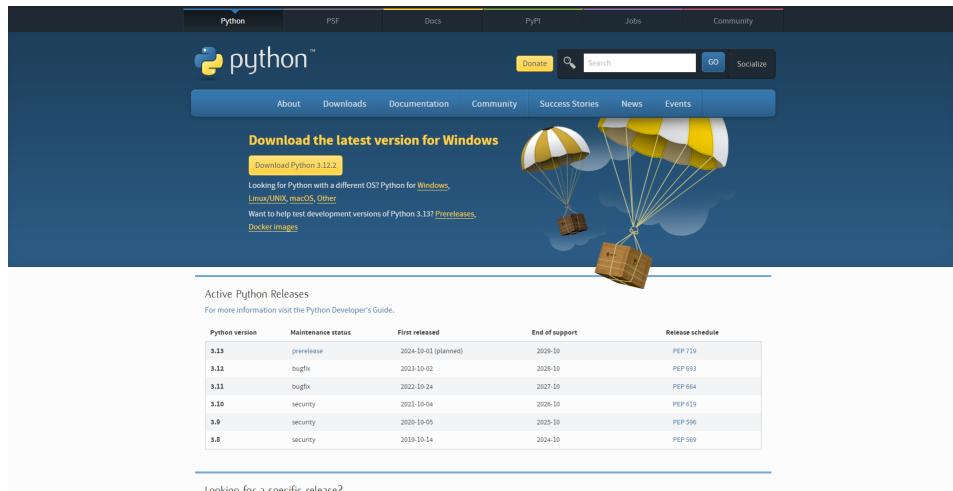


Figure 4: Download Python

Once installed, you are going to create a virtual environment, as it is good practise to avoid conflicts with other installed libraries.

Type this command in your command prompt:

```
'python3 -m venv openai-env'
```

Now, after creating the virtual environment, you need to activate it:

```
'openai-env\Scripts\activate'
```

After this step you should be able to see "openai-env" to the left of the cursor input section.

Once you have Python installed and (optionally) set up a virtual environment, the OpenAI Python library can be installed. From the command prompt, run:

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

Once this completes, running 'pip list' will show you the Python libraries you have installed in your current environment, which should confirm that the OpenAI Python library was successfully installed.

Now we are going to setup your API key. If you don't have an API key yet then you'll have to follow the instruction in API section

Open your command prompt and then insert the following command:

```
'setx OPENAI_API_KEY "your-api-key-here'
```

In order to make this key setup permanent, you ought to access Environment Variables, for that you just need to search for it in your windows search bar. Click on "New" and then set

```
"OPENAI_API_KEY"
```

as the variable name and your API key as the value.

### 3.2.2 Linux

Firstly, open your terminal and introduce the following command in order to download python:

In case you are working on Debian or Ubuntu:

```
'apt install python3 python3-dev'
```

In case you are working on Red Hat, CentOS, or Fedora:

```
'dnf install python3 python3-devel'
```

If you are having trouble installing, maybe try checking Python's beginners guide [9].

Once installed, you are going to create a virtual environment, as it is good practise to avoid conflicts with other installed libraries.

Type this command in your terminal:

```
'python3 -m venv openai-env'
```

If you can't use none of these commands above because of this error: "The virtual environment was not created successfully because ensurepip is not available" then, try using the following command (after this command you have to re-insert one of the commands above):

```
'sudo apt install python3.10-venv'
```

Now, after creating the virtual environment, you need to activate it:

```
'source openai-env/bin/activate'
```

After this step you should be able to see "openai-env" to the left of the cursor input section.

Once you have Python installed and (optionally) set up a virtual environment, the OpenAI Python library can be installed. From the terminal, run:

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

Once this completes, running 'pip list' will show you the Python libraries you have installed in your current environment, which should confirm that the OpenAI Python library was successfully installed.

Now we are going to setup your API key. If you don't have an API key yet then you'll have to follow the instruction in API section

Go to OpenAI website and access the "API keys" section, there you are going to retrieve your API key or create one in case you do not already have.

Then, open your terminal and type the following command:

```
'export OPENAI_API_KEY='your-api-key-here'',
```

To save just press Ctrl+O.

If you want to check whether it is setup correctly, type

```
'echo $OPENAI_API_KEY'
```

```

pedro@joseph:~$ python3 -m venv openai-env
pedro@joseph:~$ source openai-env/bin/activate
(openai-env) pedro@joseph:~$ pip install --upgrade openai
Collecting openai
  Downloading openai-1.14.2-py3-none-any.whl (262 kB)
262.4/262.4 KB 4.4 MB/s eta 0:00:00
Collecting sniffio
  Downloading sniffio-1.3.1-py3-none-any.whl (10 kB)
Collecting typing_extensions;>=4.7
  Downloading typing_extensions-4.10.0-py3-none-any.whl (33 kB)
Collecting aiofiles;~=3.5.0
  Downloading aiofiles-4.3.0-py3-none-any.whl (85 kB)
85.6/85.6 KB 7.7 MB/s eta 0:00:00
Collecting tqdm>4
  Downloading tqdm-4.66.2-py3-none-any.whl (78 kB)
78.3/78.3 KB 7.5 MB/s eta 0:00:00
Collecting distro<2,>=1.7.0
  Downloading distro-1.9.0-py3-none-any.whl (20 kB)
Collecting httpx<1,>=0.23.0
  Downloading httpx-0.27.0-py3-none-any.whl (75 kB)
75.0/75.0 KB 2.1 MB/s eta 0:00:00
Collecting pydantic<3,>=1.9.0
  Downloading pydantic-2.6.4-py3-none-any.whl (394 kB)
394.9/394.9 KB 6.2 MB/s eta 0:00:00
Collecting exceptiongroup>=1.0.2
  Downloading exceptiongroup-1.2.0-py3-none-any.whl (16 kB)
Collecting idna==2.8
  Downloading idna-3.6-py3-none-any.whl (61 kB)
61.6/61.6 KB 14.2 MB/s eta 0:00:00
Collecting certifi
  Downloading certifi-2024.2.2-py3-none-any.whl (163 kB)
163.8/163.8 KB 11.8 MB/s eta 0:00:00
Collecting httpcore==1.4
  Downloading httpcore-1.0.4-py3-none-any.whl (77 kB)
77.8/77.8 KB 6.2 MB/s eta 0:00:00
Collecting h11<0.15,>=0.13
  Downloading h11-0.14.0-py3-none-any.whl (58 kB)
58.3/58.3 KB 12.0 MB/s eta 0:00:00
Collecting pydantic_core==2.16.3
  Downloading pydantic_core-2.16.3-cp310-cp310-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (2.2 MB)
2.2/2.2 MB 8.5 MB/s eta 0:00:00
Collecting annotated-types<0.6.0
  Downloading annotated_types-0.6.0-py3-none-any.whl (12 kB)
Successfully installed packages: typing_extensions, tqdm, sniffio, idna, h11, exceptiongroup, distro, certifi, annotated-types, pydantic_core, httpcore, anyio, pydantic, httpx, openai, sniffio-1.3.0, tqdm-4.66.2, typing_extensions-4.10.0
(openai-env) pedro@joseph:~$ export OPENAI_API_KEY="sk-XjZZjRsV4Gg6s4lCGT3BlkFJX02L02VuYbgAxIE0EH"
(openai-env) pedro@joseph:~$ echo $OPENAI_API_KEY
sk-XjZZjRsV4Gg6s4lCGT3BlkFJX02L02VuYbgAxIE0EH
(openai-env) pedro@joseph:~$ 

```

Figure 5: Setup Python

### 3.2.3 MacOS

Firstly install Brew, if not already installed:

```
'/bin/bash -c "$(curl -fsSL https://raw.githubusercontent.com/Homebrew/install/HEAD/install.sh)"'
```

Now that you already have Brew, open your terminal and introduce the following command in order to download python:

```
'brew install python'
```

If you are having trouble installing, maybe try checking Python's beginners guide [9].

Once installed, you are going to create a virtual environment, as it is good practice to avoid conflicts with other installed libraries.

Type this command in your terminal:

```
'pip install virtualenv'
```

Create venv:

```
'virtualenv openai-env'
```

Now, after creating the virtual environment, you need to activate it:

```
'source openai-env/bin/activate'
```

After this step you should be able to see "openai-env" to the left of the cursor input section.

Once you have Python installed and (optionally) set up a virtual environment, the OpenAI Python library can be installed. From the terminal, run:

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

Once this completes, running 'pip list' will show you the Python libraries you have installed in your current environment, which should confirm that the OpenAI Python library was successfully installed.

Now we are going to setup your API key. If you don't have an API key yet then you'll have to follow the instruction in API section

Go to OpenAI website and access the "API keys" section, there you are going to retrieve your API key or create one in case you do not already have.

Then, open your terminal and type the following command:

```
'export OPENAI_API_KEY='your-api-key-here'',
```

To save just press Ctrl+O.

If you want to check it did get setup correctly, type

```
'echo $OPENAI_API_KEY'
```

### 3.3 Installing Jupyter

The tutorials in this script have been designed to be used with Jupyter [10], an interactive python environment that will ease your first steps into AI.

You can install JupyterLab, the web-based interface, by running the following command:

```
'pip install jupyterlab'
```

You can launch the interface with:

```
'jupyter lab'
```

If you are using an IDE such as *VS Code* or *Jetbrain's PyStorm* simply install the available jupyter notebook extension/plugin.

### 3.4 Installing required packages

Before starting the tutorials, install all of the necessary python packages using pip:

```
'pip install -r <path_to_requirements.txt>'
```

## 4 Tutorial/Functionality

In order to work with OpenAI's APIs you must first import it's library.

Listing 1: Import OpenAI

```
1 from openai import OpenAI
```

To use this library you need to instantiate a client. To do this you are going to need a key. If no key is indicated in the constructor, OpenAI will default to the environment variable `OPEN_AI_KEY` value.

Listing 2: Create Client

```
1 # API_KEY -> is retrieve from the environment variable
2 client = OpenAI()
3
4 # API_KEY -> is indicated by the user at the moment of creation
5 client = OpenAI(api_key="{YOUR_OPEN_AI_KEY}")
```

For good practice, we advise that you use a `.env` file to store your private information, like your key(see `.env.example`).

Listing 3: Create Client wiht .env file

```
1 import os
2 from dotenv import load_dotenv
3
4 load_dotenv()
5 api_key = os.getenv("OPENAI_API_KEY")
```

### 4.1 Chat Completions

Chat Completions is probably one of the simplest OpenAI's capabilities. It simply takes a list of messages as input and returns an answer.

#### 1 - First import and create an OpenAI client

Listing 4: Chat Completion step 1

```
1 from openai import OpenAI
2 client = OpenAI()
```

#### 2 - Structure our response, passing the gpt model and a few messages to give a brief context to our chatbot

Listing 5: Chat Completion step 2

```
1 response = client.chat.completions.create(
2     model="gpt-3.5-turbo",
3     messages=[
4         {"role": "system", "content": "You are a helpful football assistant."},
5         {"role": "user", "content": "Who won the Euro back in 2016?"},
6         {"role": "assistant", "content": "Portugal won the World Cup in 2016."},
7         {"role": "user", "content": "Where was it played, what was the score of the final and
8             who scored in that game?"}
9     ]
10 )
```

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

---

## 4.2 Assistants

The Assistants API makes it possible to create AI assistants within applications. An Assistant has instructions and can use models, tools and files to answer the user's questions. The Assistants API currently supports three types of tools: Code Interpreter, File Search and Function Calling.

### 4.2.1 Assistant

In this example, an Assistant is created as a personal maths tutor.

Initially, an Assistant is created, which represents an entity that can be configured to respond to a user's messages using various parameters such as a , 'name', 'instructions' and 'model'.

Listing 6: Assistant step 1

```
1 from openai import OpenAI
2 from runners.standard_run import std_run
3 from runners.streaming_run import streaming_run
4
5
6 client = OpenAI()
7
8 assistant = client.beta.assistants.create(
9     name="Math Tutor",
10    instructions="You are a personal math tutor. Write and run code to answer math
11      questions.",
12    model="gpt-3.5-turbo",
13 )

```

---

Next, a Thread is created that represents a conversation between a user and one or more Assistants. A Thread can be created when a user (or their AI application) starts a conversation with their Assistant.

Listing 7: Assistant step 2

```
1 thread = client.beta.threads.create()
```

---

A message is added to the Thread, the content of the messages that users or applications create is added as objects to the Message Thread.

Listing 8: Assistant step 3

```
1 message = client.beta.threads.messages.create(
2     thread_id=thread.id,
3     role="user",
4     content="I need to solve the equation: 4a - 13 = 7. Can you help me?"
5 )
```

---

When the user has added all of the necessary information/messages to the thread you must run the thread to provide them with results.

You can run a thread using its id as well as the assistant's id, like this:

Listing 9: Assistant step 4

---

```

1 run = client.beta.threads.runs.create(
2     thread_id=thread.id,
3     assistant_id=assistant.id,
4     #model = "gpt-3.5-turbo" # By default the model defined in the assistant will be
5         used. You can, however, overwrite it here.
6     instructions = "Please address the user as SISTCA student." # This parameter is
7         optional
8 )

```

---

During its life cycle, a run can have multiple states [4]:

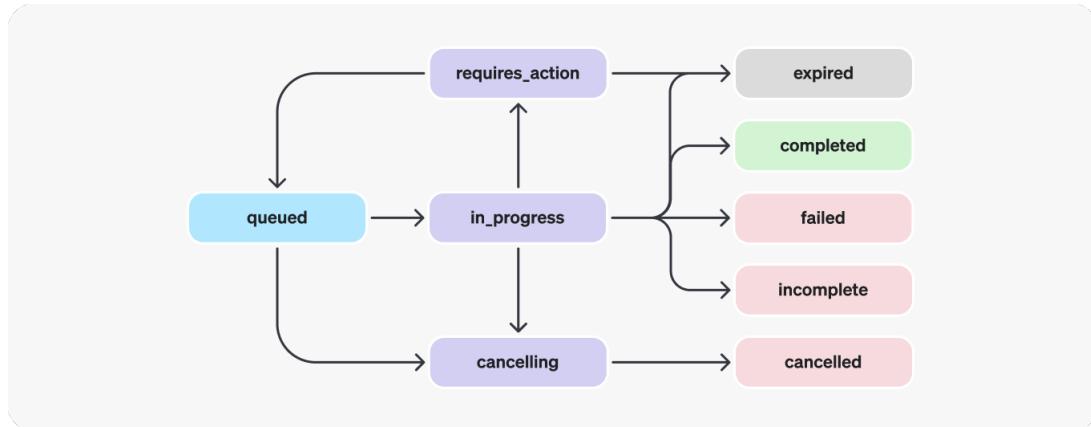


Figure 6: Run status flow diagram .

Therefore, we must check for its completion.

Listing 10: Assistant step 5

---

```

1 if run.status == 'completed':
2     messages = client.beta.threads.messages.list(
3         thread_id=thread.id
4     )
5     print(messages)
6 else:
7     print(run.status)

```

---

The result is then returned by printing the messages returned by the run.

Listing 11: Assistant step 6

---

```

1 messages = client.beta.threads.messages.list(thread_id = thread.id)
2 print(messages.data[0].content[0].text.value)

```

---

In order to yield any results, an assistant's thread must be run. This is to say that whenever you're working with the assistant functionality you must have a runner function.

Due to the fact that we will be exploring other types of assistants, and as a way of not repeating this code multiple times, we will export as the following function:

Listing 12: Assistant step 7

---

```

1 def std_run (thread_id, assistant_id, client):
2     run = client.beta.threads.runs.create_and_poll(
3         thread_id = thread_id,
4         assistant_id = assistant_id,
5         instructions = "Please address the user as SISTCA student."

```

---

```

6
7
8     if run.status == 'completed':
9         messages = client.beta.threads.messages.list(thread_id = thread_id)
10        print(messages.data[0].content[0].text.value)
11    else:
12        print(run.status)

```

---

If you have used ChatGPT before you may have noticed that the output is not returned all at once, like it is here, but rather gradually, in real time.

We can achieve this same effect by using a different type of runner, the streaming run.

In this run, the `streaming_run` function is defined, which is responsible for starting the streaming with the assistant. Within this, the '**EventHandler**' class is defined, which is inherited from **AssistantEventHandler**. This class contains the methods responsible for dealing with the different events that occur during streaming:

- '`on_text_created`': Called when a new text is created by the assistant. This prints '`assistant`' to indicate that the assistant is responding.
- '`on_text_delta`': Called for each text update. This prints the part of the text that was generated (`delta.value`).
- '`on_tool_call_created`': Called when the wizard calls a tool. This prints the type of tool call (`tool_call.type`).
- '`on_tool_call_delta`': Called for tool call updates, specifically for the `code_interpreter`. This prints the `code_interpreter` input and any outputs, especially logs.

Finally, the API client's `create_and_stream` method is used to start streaming, passing the '**EventHandler**' as the event handler.

- '`thread_id`' and '`assistant_id`' are the required identifiers.
- '`instructions`' defines specific instructions for the assistant, in this case asking it to address the user as a SISTCA student.
- '`event_handler`' is the instance of the '**EventHandler**' class.

The '`with`' block ensures that streaming continues until it is finished (`stream.until_done()`).

Listing 13: Assistant step 8

```

1 # With streaming
2
3 from typing_extensions import override
4 from openai import AssistantEventHandler
5
6
7 def streaming_run (thread_id, assistant_id, client):
8
9     # First, we create a EventHandler class to define
10    # how we want to handle the events in the response stream.
11
12    class EventHandler(AssistantEventHandler):
13        @override
14        def on_text_created(self, text) -> None:
15            print(f"\nassistant > ", end="", flush=True)
16

```

```

17     @override
18     def on_text_delta(self, delta, snapshot):
19         print(delta.value, end="", flush=True)
20
21     def on_tool_call_created(self, tool_call):
22         print(f"\nassistant > {tool_call.type}\n", flush=True)
23
24     def on_tool_call_delta(self, delta, snapshot):
25         if delta.type == 'code_interpreter':
26             if delta.code_interpreter.input:
27                 print(delta.code_interpreter.input, end="", flush=True)
28         if delta.code_interpreter.outputs:
29             print(f"\n\output >", flush=True)
30             for output in delta.code_interpreter.outputs:
31                 if output.type == "logs":
32                     print(f"\n{output.logs}", flush=True)
33
34
35     # Then, we use the 'create_and_stream' SDK helper
36     # with the 'EventHandler' class to create the Run
37     # and stream the response.
38
39     with client.beta.threads.runs.stream(
40         thread_id=thread_id,
41         assistant_id=assistant_id,
42         instructions="Please address the user as SISTCA student.",
43         event_handler=EventHandler(),
44     ) as stream:
45         stream.until_done()
46
47     print(streaming_run(thread.id, assistant.id, client))

```

---

#### 4.2.2 Function Calling

Function Calling is one of Assistants API's tools, which in simple terms, invokes predefined functions in order to respond to a certain prompt. As you will verify in just a moment, in this tutorial we asked our weather assistant to give us the maximum and minimum temperatures and rain probability in a certain location.

By using functions we allow the assistant to call additional methods as it sees fit. Assistants will analyse the user's message and decide if there is a function that can be used to answer the request and, if so, a request to call that function will be made to the runner.

For instance, in this case we want to make it so that when asking for weather data for a Portuguese city or region the data should come from IPMA.

By using IPMA'S API we can check the weather data for a specific location using it's ID.

We start by defining the functions the assistant may call, alongside their parameters.

Listing 14: Function Calling step 1

---

```

1  assistant = client.beta.assistants.create(
2      instructions="You are a weather bot. Use the provided functions to answer questions.",
3      model="gpt-3.5-turbo",
4      tools=[
5          {
6              "type": "function",
7              "function": {
8                  "name": "get_weather_data",
9                  "description": "Get the current weather forecast for the specified

```

```

10         "location",
11     "parameters": {
12         "type": "object",
13         "properties": {
14             "location": {
15                 "type": "string",
16                 "description": "City or region for which to get the weather
17                 forecast",
18             }
19         },
20         "required": ["location"]
21     }
22 }
23 )

```

---

Then, we must define actual functions to fetch the weather data.

To get weather data for a city or region we must first know its global ID. We can do this by accessing the list of identifiers and extracting the '**globalIdLocal**' for the desired city ('**location\_name**').

By visiting the list of identifiers endpoint, this *json* is returned:

Listing 15: Function Calling - IPMA's json response - step 2

```

1  {
2      "owner": "IPMA",
3      "country": "PT",
4      "data": [
5          {
6              "idRegiao": 1,
7              "idAreaAviso": "AVR",
8              "idConcelho": 5,
9              "globalIdLocal": 1010500,
10             "latitude": "40.6413",
11             "idDistrito": 1,
12             "local": "Aveiro",
13             "longitude": "-8.6535"
14         },
15         {
16             "idRegiao": 1,
17             "idAreaAviso": "BJA",
18             "idConcelho": 5,
19             "globalIdLocal": 1020500,
20             "latitude": "38.0200",
21             "idDistrito": 2,
22             "local": "Beja",
23             "longitude": "-7.8700"
24         },
25     [...]

```

---

Knowing the *json* structure, define a function to get the **globalIdLocal** for the matching city:

Listing 16: Function Calling step 3

```

1 import requests
2
3 def get_location_id(location_name):
4     url = "https://api.ipma.pt/open-data/districts-islands.json"
5     try:
6         response = requests.get(url)

```

```

7     if response.status_code == 200:
8         data = response.json()
9         for location in data["data"]:
10            # Check if the location name matches the input by converting them to
11            # lower case, so that i.e. "Porto" == "porto"
12            if location["local"].lower() == location_name.lower():
13                return str(location["globalIdLocal"])
14            return None # Location not found
15        print("Failed to fetch data:", response.status_code)
16        return None
17    except Exception as e:
18        print("An error occurred:", e)
19        return None

```

---

Test the function:

Listing 17: Function Calling step 4

```

1 location_id = get_location_id("Porto")
2 print(location_id)

```

---

The next function, 'get\_weather\_data()', uses the 'location\_id' returned by 'get\_location\_id()' to access the actual weather data of the required location, which we then return.

Listing 18: Function Calling step 5

```

1 def get_weather_data(location_id):
2     url =
3         f"https://api.ipma.pt/open-data/forecast/meteorology/cities/daily/{location_id}.json"
4     try:
5         response = requests.get(url)
6         if response.status_code == 200:
7             data = response.json()
8             forecast = data["data"][0]
9             return forecast
10        else:
11            print("Failed to fetch data:", response.status_code)
12            return None
13    except Exception as e:
14        print("An error occurred:", e)
15        return None

```

---

Test the function:

Listing 19: Function Calling step 6

```

1 weather_data = get_weather_data("1131200")
2 print(weather_data)

```

---

Having created the necessary functions we can now try running our assistant by creating a thread and adding a message requesting a weather forecast.

Listing 20: Function Calling step 7

```

1 thread = client.beta.threads.create()
2
3 #Define the location name
4 location_name = "Porto"
5
6 message = client.beta.threads.messages.create(

```

```

7     thread_id=thread.id,
8     role="user",
9     content=f"How is the weather today in {location_name}?",
10    )
11
12    print(message.content[0].text.value + "\n")

```

---

Now, initiate a run to yield a response.

Listing 21: Function Calling step 8

```

1  run = client.beta.threads.runs.create_and_poll(
2      thread_id = thread.id,
3      assistant_id = assistant.id
4  )
5
6  if run.status == 'completed':
7      messages = client.beta.threads.messages.list(
8          thread_id=thread.id
9      )
10     print(messages)
11 else:
12     print(run.status)

```

---

This time the run status returns a different value - '`required_action`'. This likely means that the assistant has called a function and is waiting for a response.

Let's check that:

Listing 22: Function Calling step 9

```

1  # Check if a run is trying to call a tool and if so, act accordingly
2  for tool in run.required_action.submit_tool_outputs.tool_calls:
3      if tool.function.name == "get_weather_data":
4
5          # let's start by checking the content of the tool call
6          print(tool.function)

```

---

The assistant wants to call the `get_weather_data` tool using the value "Porto" as the location argument.

While the assistant is smart enough to realise that it should use the `'get_weather_data'` tool to answer the user's question it does not know how to proceed, because we never linked that tool with any of our two previously defined functions.

The set of instructions to be performed when an assistant calls a function are defined in the runner. Let's do that:

Listing 23: Function Calling step 10

```

1  import json
2
3  # Check if a run is trying to call a tool and if so, act accordingly
4  for tool in run.required_action.submit_tool_outputs.tool_calls:
5      if tool.function.name == "get_weather_data":
6
7          # First, call the get_location_id function to get the location ID for the location
8          # argument in the tool call
9          location_id = get_location_id(json.loads(tool.function.arguments)[ "location" ])
10
11         # Now, knowing the location ID, we can call the get_weather_data function to get
12         # the weather data for the location

```

---

```
11     forecast = get_weather_data(location_id)
```

---

Forecast returns a *json* similar to this:

Listing 24: Function Calling step 11

```
1 {'precipitaProb': '99.0',
2  'tMin': '11.5',
3  'tMax': '16.8',
4  'predWindDir': 'W',
5  'idWeatherType': 6,
6  'classWindSpeed': 2,
7  'longitude': '-8.6294',
8  'forecastDate': '2024-05-18',
9  'classPrecInt': 2,
10 'latitude': '41.1580'}
```

---

To provide the user with relevant weather data we shall return the minimum and maximum temperature as well as the rain probability.

Listing 25: Function Calling step 12

```
1 print(f"Rain probability: {forecast['precipitaProb']}, Max temperature:
  {forecast['tMax']}, Min temperature: {forecast['tMin']}")
```

---

To output this data we first define an array in which to store the tool outputs. In this case we only have one single tool, but this would be especially useful if we were dealing with multiple tools.

Listing 26: Function Calling step 13

```
1 # Create an array to store the tool outputs
2 tool_outputs = []
```

---

Finally, we submit the output of the '`get_weather_data`' function to the run by appending an "output" key to the output array.

Listing 27: Function Calling step 14

```
1 tool_outputs.append({"tool_call_id": tool.id,
2                     "output": f"Rain probability: {forecast['precipitaProb']}, Max
  temperature: {forecast['tMax']}, Min temperature:
  {forecast['tMin']}"}
3 )
```

---

Having defined the tool outputs we must now add them to the runner so that our assistant can use them to answer our initial question.

Listing 28: Function Calling step 15

```
1 if tool_outputs:
2     try:
3         run = client.beta.threads.runs.submit_tool_outputs_and_poll(
4             thread_id=thread.id,
5             run_id=run.id,
6             tool_outputs=tool_outputs
7         )
8         print("Tool outputs submitted successfully.")
9     except Exception as e:
10         print("Failed to submit tool outputs:", e)
11 else:
```

```
12     print("No tool outputs to submit.")
```

---

We then check for the run completion and display the output message:

Listing 29: Function Calling step 16

```
1 if run.status == 'completed':
2     messages = client.beta.threads.messages.list(thread_id=thread.id)
3     print(messages.data[0].content[0].text.value)
4 else:
5     print(run.status)
```

---

The updated runner with the tool call handling code will look like this:

Listing 30: Function Calling step 17

```
1 def std_run (thread_id, assistant_id, client):
2     run = client.beta.threads.runs.create_and_poll(
3         thread_id = thread_id,
4         assistant_id = assistant_id,
5         instructions = "Please address the user as SISTCA student."
6     )
7
8     if run.status == 'completed':
9
10        messages = client.beta.threads.messages.list(
11            thread_id = thread_id
12        )
13
14        # This code acts upon the presence of image files when using code interpreter
15        if messages.data[0].content[0].type == 'image_file':
16            return(messages.data[0].content[0].image_file.file_id)
17        else:
18            print(run.status)
19
20        return(messages.data[0].content[0].text.value)
21
22    # The following code is used to handle tool/function calls
23
24    # Preparing tool outputs
25    tool_outputs = []
26    for tool in run.required_action.submit_tool_outputs.tool_calls:
27        if tool.function.name == "get_weather_data":
28            location_id = get_location_id(json.loads(tool.function.arguments)[ "location"])
29            forecast = get_weather_data(location_id)
30            tool_outputs.append({ "tool_call_id": tool.id,
31                                  "output": f"Rain probability: {forecast['precipitaProb']},
32                                             Max temperature: {forecast['tMin']}, Min temperature:
33                                             {forecast['tMax']}"})
34
35    # Submitting tool outputs and polling for completion status
36    if tool_outputs:
37        try:
38            run = client.beta.threads.runs.submit_tool_outputs_and_poll(
39                thread_id=thread_id,
40                run_id=run.id,
41                tool_outputs=tool_outputs
42            )
43            print("Tool outputs submitted successfully.")
44        except Exception as e:
45            print("Failed to submit tool outputs:", e)
```

```

45     else:
46         print("No tool outputs to submit.")
47
48 # Retrieving messages after tool output submission
49 if run.status == 'completed':
50     messages = client.beta.threads.messages.list(thread_id=thread_id)
51     print(messages.data[0].content[0].text.value)
52 else:
53     print(run.status)

```

---

Using the same logic we can also stream the response in the runner. However, since function calling uses different events that can conflict with the previously defined ones, instead of updating the '`streaming_run`' function we will create a new one.

Listing 31: Function Calling step 18

```

1 def streaming_run_fc (thread_id, assistant_id, client):
2
3     class EventHandler(AssistantEventHandler):
4         @override
5         def on_event(self, event):
6
7             # Retrieve events that are denoted with 'requires_action'
8             # since these will have our tool_calls
9
10            if event.event == 'thread.run.requires_action':
11                run_id = event.data.id # Retrieve the run ID from the event data
12                self.handle_requires_action(event.data, run_id)
13
14            def handle_requires_action(self, data, run_id):
15                tool_outputs = []
16
17                for tool in data.required_action.submit_tool_outputs.tool_calls:
18                    if tool.function.name == "get_weather_data":
19                        location_id =
20                            get_location_id(json.loads(tool.function.arguments)[ "location"])
21                        forecast = get_weather_data(location_id)
22                        tool_outputs.append({ "tool_call_id": tool.id,
23                                         "output": f"Rain probability:
24                                         {forecast['precipitaProb']}, Max temperature:
25                                         {forecast['tMin']}, Min temperature:
26                                         {forecast['tMax']}"})
27
28
29            # Submit all tool_outputs at the same time
30            self.submit_tool_outputs(tool_outputs, run_id)
31
32
33            def submit_tool_outputs(self, tool_outputs, run_id):
34                # Use the submit_tool_outputs_stream helper
35                with client.beta.threads.runs.submit_tool_outputs_stream(
36                    thread_id=self.current_run.thread_id,
37                    run_id=self.current_run.id,
38                    tool_outputs=tool_outputs,
39                    event_handler=EventHandler(),
40                ) as stream:
41                    for text in stream.text_deltas:
42                        print(text, end="", flush=True)
43                    print()
44
45
46            with client.beta.threads.runs.stream(
47                thread_id=thread_id,
48                assistant_id=assistant_id,
49            )

```

```
44     event_handler=EventHandler()
45 ) as stream:
46     stream.until_done()
```

---

#### 4.2.3 Code Interpreter

The Assistant's Code Interpreter makes it possible for assistants to write and run python code to answer user questions and analyse files.

The Code Interpreter works similarly to the tools in the previous tutorial as it too is a tool.

In this tutorial we will attempt to portray Code Interpreter's usefulness for data analysis by providing it with a file containing the data of 100 customers of a fictitious organisation.

Listing 32: Code Interpreter step 1

```
1 file = client.files.create(
2     file=open("resources/customers-100.csv", "rb"),
3     purpose='assistants'
4 )
```

---

To use code interpreter it must be passed as a tool parameter of the assistant object. We will also be adding our previously defined file as a resource.

Listing 33: Code Interpreter step 2

```
1 assistant = client.beta.assistants.create(
2     instructions="Your purpose is to analyze the provided documents and to provide
3         answers based on them.",
4     model="gpt-3.5-turbo",
5     tools=[{"type": "code_interpreter"}],
6     tool_resources={
7         "code_interpreter": {
8             "file_ids": [file.id]
9         }
10    })
```

---

Next, we create a thread and will try asking the assistant for the customer ID of Sheryl Baxter, the customer in the first row of the provided file. If it works, the value '**DD37Cf93aecA6Dc**' should be returned.

Listing 34: Code Interpreter step 3

```
1 thread = client.beta.threads.create()
2
3 message = client.beta.threads.messages.create(
4     thread_id=thread.id,
5     role="user",
6     content="What is the customer Id for Sheryl Baxter?"
7 )
8
9 print(message.content[0].text.value + "\n")
```

---

This time will simply call our previously defined runner instead of writing the runner code all over again.

Listing 35: Code Interpreter step 4

```
1 print(std_run(thread.id, assistant.id, client))
```

---

Now, for something more complex, let's ask the assistant to create a graph using the provided file.

Listing 36: Code Interpreter step 5

```
1 message = client.beta.threads.messages.create(
2     thread_id=thread.id,
3     role="user",
4     content="Create a graph of the number of customers per country."
5 )
6
7 print(message.content[0].text.value + "\n")
```

---

This action should return a file with the attached graph. When an assistant return a file it is marked as an '**image\_file**'. To access it we must get its id.

This can be achieved by checking for the presence of an '**image\_file**' in the return messages of a completed run:

Listing 37: Code Interpreter step 6

```
1 if run.status == 'completed':
2
3     messages = client.beta.threads.messages.list(
4         thread_id = thread.id
5     )
6
7     # This code acts upon the presence of image files when using code interpreter
8     if messages.data[0].content[0].type == 'image_file':
9         print(messages.data[0].content[0].image_file.file_id)
10    else:
11        print(run.status)
```

---

The updated runner with this code looks like this:

Listing 38: Code Interpreter step 7

```
1 def std_run (thread_id, assistant_id, client):
2     run = client.beta.threads.runs.create_and_poll(
3         thread_id = thread_id,
4         assistant_id = assistant_id,
5         instructions = "Please address the user as SISTCA student."
6     )
7
8     if run.status == 'completed':
9
10        messages = client.beta.threads.messages.list(
11            thread_id = thread_id
12        )
13
14        # This code acts upon the presence of image files when using code interpreter
15        if messages.data[0].content[0].type == 'image_file':
16            return(messages.data[0].content[0].image_file.file_id)
17        else:
18            print(run.status)
19
20        return(messages.data[0].content[0].text.value)
21
22    # The following code is used to handle tool/function calls
23
```

```

24     # Preparing tool outputs
25     tool_outputs = []
26     for tool in run.required_action.submit_tool_outputs.tool_calls:
27         if tool.function.name == "get_weather_data":
28             location_id = get_location_id(json.loads(tool.function.arguments)[ "location" ])
29             forecast = get_weather_data(location_id)
30             tool_outputs.append({ "tool_call_id": tool.id,
31                                 "output": f"Rain probability: {forecast['precipitaProb']},"
32                                         "Max temperature: {forecast['tMin']}, Min temperature:"
33                                         f"{forecast['tMax']}"} )
34
35     # Submitting tool outputs and polling for completion status
36     if tool_outputs:
37         try:
38             run = client.beta.threads.runs.submit_tool_outputs_and_poll(
39                 thread_id=thread_id,
40                 run_id=run.id,
41                 tool_outputs=tool_outputs
42             )
43             print("Tool outputs submitted successfully.")
44         except Exception as e:
45             print(f"Failed to submit tool outputs: {e}")
46     else:
47         print("No tool outputs to submit.")
48
49     # Retrieving messages after tool output submission
50     if run.status == 'completed':
51         messages = client.beta.threads.messages.list(thread_id=thread_id)
52         print(messages.data[0].content[0].text.value)
53     else:
54         print(run.status)

```

---

Store the file ID:

Listing 39: Code Interpreter step 8

```
1 file_id = std_run(thread.id, assistant.id, client)
```

---

To retrieve the file we must call the '`files_endpoint`' using the '`file_id`'.

Listing 40: Code Interpreter step 9

```
1 image_data = client.files.content(file_id=file_id)
2 image_data_bytes = image_data.read()
```

---

Save and display the result:

Listing 41: Code Interpreter step 10

```

1 import IPython
2
3 image_name = "graph.png"
4
5 image_path = os.path.join("resources", image_name)
6
7 with open(image_path, "wb") as image_file:
8     image_file.write(image_data_bytes)
9
10 IPython.display.Image(image_name)

```

---

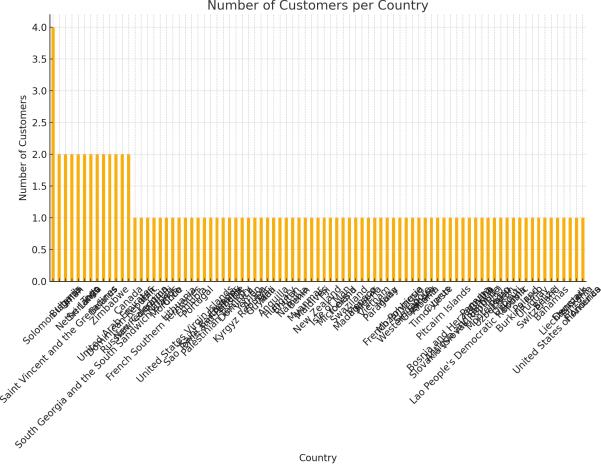


Figure 7: Resulting graph

#### 4.2.4 File Search

File search is a tool that complements the Assistants API. It allows users to send their documents and question the assistant about them. The documents sent by the user are transformed in embeddings and stored in vector stores. File search uses both keyword and vector search to retrieve relevant information.

This feature implements good retrieval practices to extract the necessary data to answer the user query:

- Rewrites user queries to optimise them for search.
- Breaks down complex user queries into multiple searches that can run in parallel.
- Runs both keyword and semantic searches across both assistant and thread vector stores.
- Reranks search results to pick the most relevant ones before generating the final response.

File search supports a variety of different file formats. For this tutorial, a simple txt containing some wikipedia articles was the source of the information.

- .c
- .cs
- .cpp
- .doc
- .docx
- .html
- .java
- .json
- .md
- .pdf
- .php

- .pptx
- .py
- .rb
- .tex
- .txt
- .css
- .js
- .sh
- .ts

## Vector Store

A vector store is a digital storage system that holds vectors. Adding a file to a vector store automatically parses, chunks, embeds and stores the file in a vector database that is capable of both keyword and semantic search. Adding a file or files to a vector store is an *async* operation. The maximum file size you can store 512 MB and each file can not contain more than 5,000,000 tokens, this is computed automatically when a file is attached.

In this tutorial we will provide the assistant with a text file containing wikipedia articles about musical artists. We will then have the assistant search this file to answer questions.

Listing 42: File Search step 1

---

```

1  assistant = client.beta.assistants.create(
2      name="Music Information Assistant",
3      instructions="You are an expert music analyst. Use your knowledge base to answer
4          questions regarding musical artists and their work.",
5      model="gpt-3.5-turbo",
6      tools=[{"type": "file_search"}],
7

```

---

Create a vector store to store your information

Listing 43: File Search step 2

---

```

1  vector_store = client.beta.vector_stores.create(name="Music Information")
2
3  # Add the desired files. It is possible to add multiple files
4  file_paths = ["resources/file_search_information.txt"]
5  file_streams = [open(path, "rb") for path in file_paths]
6
7
8  # upload the files and add them to the vector store
9  file_batch = client.beta.vector_stores.file_batches.upload_and_poll(
10     vector_store_id=vector_store.id, files=file_streams
11 )

```

---

Then, update the assistant with the newly created vector store.

Listing 44: File Search step 3

---

```

1  assistant = client.beta.assistants.update(
2      assistant_id=assistant.id,
3      tool_resources={"file_search": {"vector_store_ids": [vector_store.id]}},
4

```

```

5
6 # upload the user provided file to OpenAI
7 message_file = client.files.create(
8     file=open("resources/file_search_information.txt", "rb"), purpose="assistants"
9 )

```

---

Create a thread and ask the assistant a question whilst attaching the file to the input.

Listing 45: File Search step 4

```

1 thread = client.beta.threads.create(
2     messages=[
3         {
4             "role": "user",
5             "content": "How many artists do you know about, and what are their albums?",
6             "attachments": [
7                 { "file_id": message_file.id, "tools": [{"type": "file_search"}] }
8             ],
9         }
10    ]
11 )

```

---

Lastly, run the thread and print the response:

Listing 46: File Search step 5

```

1 print(std_run(thread.id, assistant.id, client))

```

---

### 4.3 Embeddings

OpenAI's text embeddings transforms strings to numbers, that allows to measure the relatedness of text strings. Embeddings are commonly used for:

- Search (where results are ranked by relevance to a query string)
- Clustering (where text strings are grouped by similarity)
- Recommendations (where items with related text strings are recommended)
- Anomaly detection (where outliers with little relatedness are identified)
- Diversity measurement (where similarity distributions are analyzed)
- Classification (where text strings are classified by their most similar label)

An embedding is a vector (list) of floating point numbers. The distance between two vectors measures their relatedness. Small distances suggest high relatedness and large distances suggest low relatedness.

Like mentioned before, Embeddings have a lot of uses but this tutorial will only focus on how to make the request and a show simple comparison between two strings.

#### Models

Right now there are three available models for Embeddings. The ones with "-3" on the name are third generation models.

- text-embedding-3-small

- text-embedding-3-large
- text-embedding-ada-002

## 1 - Import the necessary libraries and start the client

Listing 47: Embeddings tutorial step 1

---

```
1 from openai import OpenAI
2 client = OpenAI()
```

---

## 2 - Retrieve the embeddings

*model* - changes the model use to retrieve information *input* - string message you want to retrieve the embeddings from

Listing 48: Embeddings tutorial step 2

---

```
1 response1 = client.embeddings.create(
2     input="We are testing to see if this string has any similarities to another one.",
3     model="text-embedding-3-small"
4 )
5
6 embeddings1 = response1.data[0].embedding
7
8 print(f"Embeddings Example - {embeddings1[0]}, {embeddings1[1]}")
```

---

## Compare two strings using embeddings and cosine similarity

### 1 - Import numpy and cosine\_similarity

Listing 49: Embeddings Compare Two Strings step 1

---

```
1 import numpy as np
2 from sklearn.metrics.pairwise import cosine_similarity
```

---

### 2 - Retrieve the embeddings for the second string

Listing 50: Embeddings Compare Two Strings step 2

---

```
1 response2 = client.embeddings.create(
2     input="We're experimenting to determine if this string bears resemblance to another.",
3     model="text-embedding-3-small"
4 )
5
6 embeddings2 = response2.data[0].embedding
```

---

### 3 - Convert the embeddings to numpy arrays

Listing 51: Embeddings Compare Two Strings step 3

---

```
1 embeddings1 = np.array(embeddings1).reshape(1, -1)
2 embeddings2 = np.array(embeddings2).reshape(1, -1)
```

---

### 4 - Calculate the similarity

The closer the score is to 1 the similar it is.

Listing 52: Embeddings Compare Two Strings step 4

---

```
1 similarity_score = cosine_similarity(embeddings1, embeddings2)
2
3 final_score = float(format(similarity_score[0][0], ".2f"))
4
5 print(f"Similarity: {final_score}")
```

---

## 4.4 Vision

Vision API from OpenAI, based on GPT-4 Turbo, is an enormous jump in AI capability. This functionality brought AI the capability to take pictures as input, comprehend them and answer questions about them. This development guarantees to revolutionize the way we interact and use AI.

Vision can accept images through links or even by passing its base64 encoded image. Unfortunately, this feature is only available to ChatGPT-4, so we will be addressing this tutorial at the end of our article as an extra.

1. Firstly, in order to start using this functionality, you need to import the necessary packages. In this case, we only need to import the '**openai**'.

---

```
1 from openai import OpenAI
```

---

2. Now it's time to create a '**client**' in order to request the '**response**' from '**openai**'.

---

```
1 client = OpenAI()
```

---

3. Last but not least, the '**response**' is structured, passing the role, content and the image.

---

```
1 response = client.chat.completions.create(
2     model="gpt-4-turbo",
3     messages=[
4         {
5             "role": "user",
6             "content": [
7                 {"type": "text", "text": "What's in this image?"},
8                 {
9                     "type": "image_url",
10                    "image_url": {
11                        "url":
12                            "https://upload.wikimedia.org/wikipedia/commons/thumb/d/dd/Gfp-wisconsin-madison-"
13                                },
14                            },
15                        ],
16                    },
17                    max_tokens=300,
18                )
19
20    print(response.choices[0].message.content)
```

---

## 4.5 Image Generation (DALL-E)

The open AI image API provides three methods for interacting with images, namely creating images from scratch using a text prompt (DALL-E 3 and DALL-E 2), creating edited versions of

images by having the model change some areas of a pre-existing image based on a new text prompt (DALL-E 2 only) and creating variations of an existing image (DALL-E 2 only).

This guide covers the basics of the API methods with useful code examples.

#### 4.5.1 Image Generation

The image generation method allows you to create an original image with a text prompt.

Initially, it is necessary to import the display and Image functions from the IPython.display module, in order to display the image inside a Jupyter notebook.

It also imports os, used for operating system operations, as explained earlier for creating a venv environment, and imports the OpenAI class from the openai module, an interface for using the API.

The code creates an instance of the OpenAI() client. A request is made to the API to generate an image, with the parameters:

- **model:** specifies the model used to generate the image, in this case "dall-e-2", which specialises in generating images based on textual descriptions;
- **prompt:** descriptive text that will serve as input for the model to generate the image, in this example, "A cat inside a car";
- **size:** specifies the size, 1024x1024 pixels;
- **quality:** sets the quality, "standard". When using DALL-E 3 it is possible to set quality: "hd", i.e. fine detail. However, standard quality square images are generated more quickly;
- **n:** defines the number of images generated.

After the API generates the image, the image URL is extracted from the response and stored in the *image\_url* variable.

Finally, the image URL is printed and the image is displayed using the display function with the Image class, passing the URL as a parameter and setting the image width to 500 pixels.

---

```
1  from IPython.display import display, Image
2
3  from openai import OpenAI
4
5  client = OpenAI()
6
7  response = client.images.generate(
8      model="dall-e-2",
9      prompt="A cat inside a car",
10     size="1024x1024",
11     quality="standard",
12     n=1,
13 )
14
15 image_url = response.data[0].url
16 print(image_url)
17 display(Image(url=image_url, width=400))
```

---

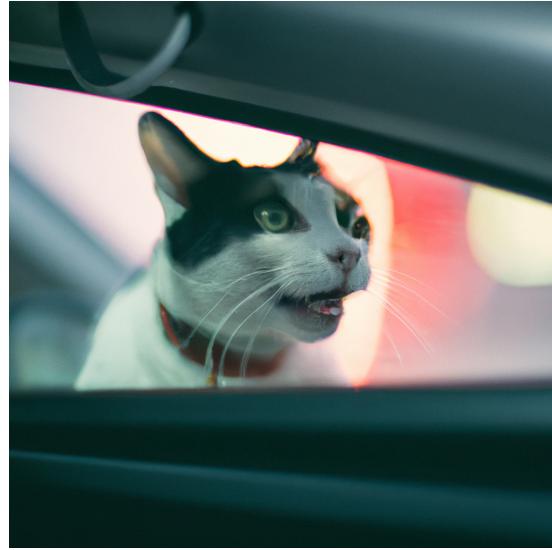


Figure 8: Image generated

#### 4.5.2 Edit (DALLE 2 only)

Also known as "inpainting", the image editing endpoint allows you to edit an image by loading an image and a mask indicating the areas that should be replaced, you can use tools such as [GIMP] or , to erase a certain area of the image. The transparent areas of the mask indicate where the image should be edited, and the prompt should describe the complete new image, not just the deleted area.

The image and mask sent must be square PNG images of less than 4 MB and must also have the same dimensions. The non-transparent areas of the mask are not used to generate the output, so they don't necessarily have to match the original image as in the following example, which shows the original image, the mask and the resulting image after the editing process.



Figure 9: Input/Output from the edit endpoint

---

```

1  from openai import OpenAI
2  client = OpenAI()
3  from IPython.display import display, Image
4
5  response = client.images.edit(
6      model="dall-e-2",
7      image=open("original.png", "rb"),
8      mask=open("mask.png", "rb"),
9      prompt="A hand holding a sandwich",
10     n=1,
11     size="1024x1024"

```

```

12     )
13     image_url = response.data[0].url
14
15     print(image_url)
16     display(Image(url=image_url, width=400))

```

---

#### 4.5.3 Variations (DALL·E 2 only)

The image variations endpoint allows you to generate a variation of a given image.

Similar to the edits endpoint, the input image must be a square PNG image less than 4MB in size.

```

1  from IPython.display import display, Image
2  from openai import OpenAI
3  client = OpenAI()
4
5  response = client.images.create_variation(
6      model="dall-e-2",
7      image=open("original.png", "rb"),
8      n=1,
9      size="1024x1024"
10 )
11
12 image_url = response.data[0].url
13
14 print(image_url)
15 display(Image(url=image_url, width=400))

```

---

The images below correspond to a possible example of the endpoint variations.



Figure 10: Input/Output from the variation endpoint

## 4.6 TTS

In this tutorial you will learn how to integrate text-to-speech feature from open-ai in your projects.

Voice only changes the tone and the "person" who is speaking. TTS only produces english audio files.

There are two models available tts right now: **tts-1** and **tts-1-hd**. If you want lower latency **tts-1** is recommended, but it comes with lower quality than **tts-1-hd**.

Available voices	Output Formats
Alloy	<b>mp3</b> - default format
Echo	<b>opus</b>
Fable	<b>aac</b>
Onyx	<b>flac</b>
Nova	<b>wav</b>
Shimmer	<b>pcm</b>

Table 3: TTS Information Table

### 1- Importing the necessary libraries and creating the client

*pathlib* offers classes that represent a filesystem path.

Listing 53: TTS tutorial step 1

---

```

1 from pathlib import Path
2 from openai import OpenAI
3
4 client = OpenAI()

```

---

### 2 - Create a path to save the audio file

In here you can choose the format you want your audio file to be in.

Listing 54: TTS tutorial step 2

---

```

1 speech_file_path = Path(f"Tutorials/TTS/tts_audio.mp3")

```

---

### 3 - Generate the audio

The endpoint will receive a model, a voice and your input.

Listing 55: TTS tutorial step 3

---

```

1 response = client.audio.speech.create(
2     model="tts-1",
3     voice="shimmer",
4     input="Hey, I'm a student in Lincenciatura de Engenharia de Telecomunicacoes e
      Informatica in Instituto Superior de Engenharia do Porto, and I'm doing a tutorial
      on how to use open-a.i in my projects!"
5 )
6
7 response.write_to_file(speech_file_path)

```

---

## 4.7 Whisper

In this tutorial you will learn how to use Whisper to transcribe text from audio files as well translate it into English.

---

```

1 import os
2 from openai import OpenAI
3
4 client = OpenAI()

```

---

### Load the audio file

Feel free to try different audio files as well as add/record your own. Files must be of one of these types: mp3, mp4, mpeg, mpg, m4a, wav, and webm.

**Note that files greater than 25MB will need to be segmented using additional libraries**

```
1     audio_file = open("audio.wav", "rb")
```

## 1. Transcribe an audio file

The transcription endpoint will take the input audio and transcribe it into text.

```
1     transcription = client.audio.transcriptions.create(
2         model="whisper-1",
3         file=audio_file,
4         response_format="text"
5     )
6
7     print(transcription)
```

Notice how `response_format="text"`? To get additional information to get additional information try changing it to `verbose_json`.

You should now receive a json response with additional parameters. One of which, the `language` parameter, includes the detected language from the input file.

**Note:** If the language is not being properly detected, which may negatively impact transcription, you can add an additional parameter stating it according to the ISO-639-1 format.

```
1     (
2         model="whisper-1",
3         file=audio_file,
4         response_format="text"
5         language="..."
6     )
```

We can modify our code to reflect this:

```
1     transcription = client.audio.transcriptions.create(
2         model="whisper-1",
3         file=audio_file,
4         response_format="verbose_json"
5     )
6     print(f"Detected language: {transcription.language}")
7     print(transcription.text)
```

## 2. Translation

Using the translation endpoint we can translate the contents of the audio file to English (currently this the only available language for translation).

```
1     translation = client.audio.translations.create(
2         model = "whisper-1",
3         file = audio_file,
4         response_format="text"
5     )
6     print(translation)
```

## 4.8 Moderation

The moderation endpoint is a tool you can use to check whether text is potentially harmful. It can be used to identify content that could be harmful and take action.

The templates classify the following categories:

- **Hate** - Content that expresses or promotes hatred based on race, gender, ethnicity, religion, nationality, sexual orientation, disability status or caste. Hateful content directed at unprotected groups constitutes harassment.
- **Hate/Threatening** - Hateful content that also includes violence or serious harm towards the targeted group based on race, gender, ethnicity, religion, nationality, sexual orientation, disability status, or caste.
- **Harassment** - Content that expresses, incites, or promotes harassing language towards any target.
- **Harassment/Threatening** - Harassment content that also includes violence or serious harm towards any target.
- **Self-harm** - Content that promotes, encourages, or depicts acts of self-harm, such as suicide, cutting, and eating disorders.
- **Self-harm/Intent** - Content where the speaker expresses that they are engaging or intend to engage in acts of self-harm, such as suicide, cutting, and eating disorders.
- **Self-harm/Instructions** - Content that encourages performing acts of self-harm, such as suicide, cutting, and eating disorders, or that gives instructions or advice on how to commit such acts.
- **Sexual** - Content meant to arouse sexual excitement, such as the description of sexual activity, or that promotes sexual services (excluding sex education and wellness).
- **Sexual/Minors** - Sexual content that includes an individual who is under 18 years old.
- **Violence** - Content that depicts death, violence, or physical injury.
- **Violence/Graphic** - Content that depicts death, violence, or physical injury in graphic detail.

To obtain a classification for a piece of text, a request is made to the moderation endpoint, as shown in the following code fragment. Firstly, the OpenAI class from the openai module is imported, then an instance of the class is created and assigned to a variable, configuring the client to interact with the OpenAI API.

With the '**moderations.create()**' method, the text is sent for moderation.

The response from the API will be a JSON object with information about the moderation of the text.

```
1  from openai import OpenAI
2
3  client = OpenAI()
4
5  response = client.moderations.create(
6      input="I own a lot of guns which I intend to use, in order to harm people"
7  )
8
9  output = response.results[0]
10 print(output)
```

Below is an example of the endpoint response structure. It returns the following fields:

- **flagged**: Set to true if the model classifies the content as potentially harmful, false otherwise.
- **categories**: Contains a dictionary of violation flags for each category. The value will be true if the model flags the corresponding category as violated, false otherwise.
- **category\_scores**: Contains a dictionary of raw scores per category issued by the model, denoting the model's confidence that the input violates the OpenAI policy for the category. The value is between 0 and 1, with higher values denoting greater confidence. The scores should not be interpreted as probabilities.

---

```
1  {
2      "id": "modr-XXXXX",
3      "model": "text-moderation-007",
4      "results": [
5          {
6              "flagged": true,
7              "categories": {
8                  "sexual": false,
9                  "hate": false,
10                 "harassment": false,
11                 "self-harm": false,
12                 "sexual/minors": false,
13                 "hate/threatening": false,
14                 "violence/graphic": false,
15                 "self-harm/intent": false,
16                 "self-harm/instructions": false,
17                 "harassment/threatening": true,
18                 "violence": true
19             },
20             "category_scores": {
21                 "sexual": 1.2282071e-6,
22                 "hate": 0.010696256,
23                 "harassment": 0.29842457,
24                 "self-harm": 1.5236925e-8,
25                 "sexual/minors": 5.7246268e-8,
26                 "hate/threatening": 0.0060676364,
27                 "violence/graphic": 4.435014e-6,
28                 "self-harm/intent": 8.098441e-10,
29                 "self-harm/instructions": 2.8498655e-11,
30                 "harassment/threatening": 0.63055265,
31                 "violence": 0.99011886
32             }
33         }
34     ]
35 }
36 }
```

---

## 5 Exercises

### 5.1 Exercise A

Embeddings have a lot of uses, and when combined with other APIs, they can do even more. One example is using embeddings with chat completion to extract information from a PDF and then creating a function to ask anything about the document.

In the following exercise, you will create a program that retrieves information from a PDF and answers questions about it. In order to achieve this, you must:

- Convert a PDF file to embeddings and save them in a CSV file.
- Use embeddings to search a user query in the CSV file.
- Send that information to chat completions.

The PDF used in this exercise will be `LETI_SISTCA_2023_24_Team2_OpenAI.pdf`.

#### 1. Start by importing the requiring dependecies and initialize the client and creating some constants

Listing 56: Exercise A - Part 1 Step 1

```
1  from openai import OpenAI
2  import pandas as pd
3  import re
4  import tiktoken
5  import PyPDF2
6
7  import ast
8  from scipy import spatial
9
10 client = OpenAI()
11
12
13 SECTIONS_TO_IGNORE = [
14     "Contents",
15     "List of Tables",
16     "List of Figures",
17     "References",
18 ]
19
20 MAX_TOKENS = 1600
21 BATCH_SIZE = 1000
22
23 # TODO #1: Create consts for GPT_MODEL and EMBEDDING_MODEL (small)
```

#### 2. Simple Logic to Extract the information from the pdf

This is a simple logic to extract the necessary information from the pdf we are going to use.

Listing 57: Exercise A - Part 1 Step 2

```
1 def extract_text_from_pdf(pdf_path):
2     with open(pdf_path, 'rb') as file:
3         reader = PyPDF2.PdfReader(file)
4         text = ''
5         for page_num in range(len(reader.pages)):
6             text += reader.pages[page_num].extract_text()
7     return text
```

```

8
9     def split_sections_from_pdf(pdf_text):
10    title = []
11    text = []
12    ignore = True
13    current_section = "I.N.I.T.I.A.L-V.A.L.U.E"
14    for line in pdf_text.split('\n'):
15        line = line.strip()
16
17        if not line:
18            continue
19        if is_new_section(line):
20            ignore = ignore_section(line)
21            if ignore:
22                continue
23            title.append(line)
24            if current_section != "I.N.I.T.I.A.L-V.A.L.U.E":
25                text.append(current_section)
26            current_section = ""
27        else:
28            if not ignore:
29                current_section += " " + line
30    if current_section:
31        text.append(current_section)
32
33
34    sections = [(title),(text)]
35    return sections
36
37    def ignore_section(line):
38        if any(section in line for section in SECTIONS_TO_IGNORE):
39            return True
40        return False
41
42    def is_new_section(line):
43        pattern = r"\d+\.\d+(?:\.\d+)? [A-Z].*?"
44
45        if line.strip().count(',') > 7:
46            return False
47        if re.match(pattern, line.strip()):
48            return True
49        if any(section in line for section in SECTIONS_TO_IGNORE):
50            return True
51
52        return False
53
54    def clean_section(section):
55
56        titles = section[0]
57        text = section[1]
58
59        for line in text:
60            line = re.sub(r"\[\d+\]", "", line)
61            line = re.sub(r"\[\d\d+\]", "", line)
62            line = line.strip()
63
64        return (titles, text)
65
66
67    def num_tokens(text, model = GPT_MODEL):
68        encoding = tiktoken.encoding_for_model(model)
69        return len(encoding.encode(text))
70

```

```

71     def halved_by_delimiter(string, delimiter = "\n"):
72         chunks = string.split(delimiter)
73         if len(chunks) == 1:
74             return [string, ""]
75         elif len(chunks) == 2:
76             return chunks
77         else:
78             total_tokens = num_tokens(string)
79             halfway = total_tokens // 2
80             best_diff = halfway
81             for i, chunk in enumerate(chunks):
82                 left = delimiter.join(chunks[: i + 1])
83                 left_tokens = num_tokens(left)
84                 diff = abs(halfway - left_tokens)
85                 if diff >= best_diff:
86                     break
87                 else:
88                     best_diff = diff
89             left = delimiter.join(chunks[:i])
90             right = delimiter.join(chunks[i:])
91             return [left, right]
92
93
94     def truncated_string(string, model, max_tokens, print_warning = True,):
95         encoding = tiktoken.encoding_for_model(model)
96         encoded_string = encoding.encode(string)
97         truncated_string = encoding.decode(encoded_string[:max_tokens])
98         if print_warning and len(encoded_string) > max_tokens:
99             print(f"Warning: Truncated string from {len(encoded_string)} tokens to
{max_tokens} tokens.")
100        return truncated_string
101
102
103    def split_strings_from_subsection(title, text, max_tokens = 1000, model =
GPT_MODEL, max_recursion = 5):
104
105        string = "\n\n".join(title + text)
106        num_tokens_in_string = num_tokens(string, model)
107
108        if num_tokens_in_string <= max_tokens:
109            return [string]
110
111        elif max_recursion == 0:
112            return [truncated_string(string, model, max_tokens)]
113
114        else:
115            for delimiter in ["\n\n", "\n", ". "]:
116                left, right = halved_by_delimiter(text, delimiter=delimiter)
117                if left == "" or right == "":
118
119                    continue
120                else:
121
122                    results = []
123                    for half in [left, right]:
124                        half_strings = split_strings_from_subsection(title,
half,max_tokens,model,max_recursion - 1)
125                        results.extend(half_strings)
126                    return results
127
128    return [truncated_string(string, model, max_tokens)]

```

---

3. Call all functions to retrieve the clean pdf sections to then split it into strings

---

Listing 58: Exercise A - Part 1 Step 3

---

```
1 # TODO #2: Call the previous created functions
2 pdf_text = ...
3 pdf_sections = ...
4 cleaned_sections = ...
5
6 MAX_TOKENS = 1600
7 strings = []
8 titles = cleaned_sections[0]
9 texts = cleaned_sections[1]
10 for i in range(len(titles)):
11     strings.extend(split_strings_from_subsection(titles[i], texts[i],
12                                                 max_tokens=MAX_TOKENS))
```

---

4. Transforming the information to embeddings and saving it to a CSV file

---

Listing 59: Exercise A - Part 1 Step 4

---

```
1
2
3     embeddings = []
4     for batch_start in range(0, len(strings), BATCH_SIZE):
5         batch_end = batch_start + BATCH_SIZE
6         batch = strings[batch_start:batch_end]
7
8     # TODO #3: Make a request to the embeddings API with the batch as input
9
10    for i, be in enumerate(response.data):
11        assert i == be.index
12    batch_embeddings = [e.embedding for e in response.data]
13    embeddings.extend(batch_embeddings)
14
15
16    df = pd.DataFrame({"text": strings, "embedding": embeddings})
17
18
19    SAVE_PATH = "SISTCA_TEAM2.csv"
20    df.to_csv(SAVE_PATH, index=False)
```

---

5. The first part is done, now we need to create a function so GPT can awnser anything about the pdf using the saved embeddings

Change the Embedding Model and read the CSV file

---

Listing 60: Exercise A - Part 2 Step 1

---

```
1 # TODO #4: Change the EMBEDDING_MODEL (ada)
2
3
4 # TODO #5: Create a variable with the CSV file path
5
6
7 df = pd.read_csv(embeddings_path)
8 df['embedding'] = df['embedding'].apply(ast.literal_eval)
```

---

6. Functions to compare the relatedness off the strings with the query

---

Listing 61: Exercise A - Part 2 Step 2

---

```

2     def strings_ranked_by_relatedness(query, df, relatedness_fn = lambda x, y: 1 -
3         spatial.distance.cosine(x, y), top_n = 100) :
4
5
6
7     query_embedding = query_embedding_response.data[0].embedding
8     strings_and_relatednesses = [
9         (row["text"], relatedness_fn(query_embedding, row["embedding"]))
10        for i, row in df.iterrows()
11    ]
12    strings_and_relatednesses.sort(key=lambda x: x[1], reverse=True)
13    strings, relatednesses = zip(*strings_and_relatednesses)
14    return strings[:top_n], relatednesses[:top_n]
15
16 strings, relatednesses = strings_ranked_by_relatedness("open ai", df, top_n=5)
17
18
19 def num_tokens(text, model = GPT_MODEL):
20     encoding = tiktoken.encoding_for_model(model)
21     return len(encoding.encode(text))
22
23 def query_message(query, df, model, token_budget):
24     strings, relatednesses = strings_ranked_by_relatedness(query, df)
25     introduction = 'Use the below articles on the document about OpenAI made by
26     Team 2, composed by Patricia Sousa, Carlos Alves, Jose Leal and Tiago
27     Ribeiro, for SISTCA to answer the subsequent question. If the answer
28     cannot be found in the articles, write "Sorry, the information you seek
29     cannot be found in the document in question."
30     question = f"\n\nQuestion: {query}"
31     message = introduction
32     for string in strings:
33         next_article = f'\n\nPDF article section:\n"""\n{string}\n"""\n'
34         if (
35             num_tokens(message + next_article + question, model)
36             > token_budget
37         ):
38             break
39         else:
40             message += next_article
41     return message + question

```

---

## 7. Create the final function

If the chat completions do not have information to answer the query, as defined, it will say “Sorry, the information you seek cannot be found in the document in question.”

Listing 62: Exercise A - Part 2 Step 3

---

```

1     def ask(query, df = df, model = GPT_MODEL, token_budget = 4096 - 500,
2            print_message = False):
3         message = query_message(query, df, model=model, token_budget=token_budget)
4         if print_message:
5             print(message)
6         messages = [
7             {"role": "system", "content": "You answer questions about the document
8              made by Team2 for SISTCA about OpenAI."},
9             {"role": "user", "content": message},
10            ]
11
12 # TODO 7: Make a request to the Chat Completions API
13

```

```
14     response_message = response.choices[0].message.content
15     return response_message
```

---

#### 8. Test the ask function to verify it

Because every time you run ask you give a new prompt to chat completions the answers may vary for one attempt to another.

Listing 63: Exercise A - Part 2 Step 4

---

```
1  print(ask("Scientific/technological background"))
2  print(ask("Give me the authors"))
3  print(ask("What can you tell me about the document"))
4  print(ask("Give me the document structure"))
```

---

The solution for this exercise is available here on [github](#).

## 5.2 Exercise B

Despite, in it's current state, only allowing English translation, we can increment the Whisper API with other APIs so as to be able to translate audio into other languages. By using the standard GPT model in the Chat Completions API we can translate text to and from any language.

In the following exercise you will create a program that translates an audio message into any other language. In order to achieve this you must:

- Use Whisper to transcribe the original audio;
- Translate the resulting transcription into any language using Chat Completions;
- Return the translated audio using Text-To-Speech.

#### 1. Start by importing the required dependencies and initializing the client

---

```
1  import os
2  from dotenv import load_dotenv
3  from openai import OpenAI
4
5  load_dotenv()
6
7  api_key = os.getenv("OPENAI_API_KEY")
8
9  client = OpenAI(api_key=api_key)
```

---

#### 2. Use Whisper to transcribe the original audio

---

```
1  # TODO #1: load the audio file
2
3
4  # TODO #2: transcribe contents and detect the input language
```

---

#### 3. Translate the transcribed text into any language using Chat Completions

---

```
1  # TODO #3: translate the transcription to another language using
2  #             chatCompletions
```

---

#### 4. Return the translated audio using Text-To-Speech

---

```
1      # TODO #4: output the translated audio file
2
3
4
5
6      # Save the output audio file
7      output_file_path = Path(__file__).parent / f"translation.mp3"
8      translated_audio.write_to_file(output_file_path)
9
10
11
12      # Test the output
13      import IPython.display as display
14
15      print(f"Detected language: {detected_language}\n")
16      print(f"Transcribed text: {transcribed_text}\n")
17
18      print(f"Translated text: {translated_text}\n")
19      print(f"Translated audio file saved to {output_file_path}")
20
21
22      display.Audio(speech_file_path)
```

---

The solution for this exercise can be found [here](#).

## 6 Challenge

With all tutorials done, you should now be able to answer correctly to this final challenge. In this section we propose you the following exercise: You have to create an assistant which will be capable of analysing and answering questions about films, series and books. The average question should look like this, 'Give me the top 10 best films of all time.'

We advise that you follow this steps:

- Create an Assistant with function calling.
  - For a better flexibility you should retrieve from the user query the type (movie, series, book), theme(comedy, thriller...) and quantity.
- Create a function to retrieve information from Chat Completions.
  - You can force Chat Completions to send the information in a specific format, like JSON.
- Show the response in a bullet list.

In our github you can find a Flask template where you can visually test the results of your challenge. To have a compatibility with the template, make sure that your final function, the one that retrieves the assistant message, follows this format:

Listing 64: Challenge Flask Function

---

```
1 def get_cinema_info(type_filter, genre_filter, quantity_filter):  
2     ...
```

---

## 7 Future Work

The OpenAI API is still an in development effort this means that everything is prone to change. In fact, some of it has changed during the development of this script. As previously noted before some of the functionalities that have been addressed is not yet available for the explorer tier of the API, which includes the "GPT-3.5-Turbo" model. At the time of writing this, the launch of a new flagship model is on the horizon - "GPT-4o" [11]. This should make previously tier-locked APIs such as Vision accessible to everyone, as well as potentially added new features.

Despite the rather wide scope of this script some of OpenAI's tools haven't been explored due to them being unavailable in the API. Such is the case for Sora [12], a video generation model, or the research publications like CLIP [13, 14], Jukebox [15, 16] or Point-E [17, 18]. Another feature unable to be explored, despite being compatible with "GPT-3.5-Turbo", was fine-tuning. This feature allows to train "GPT" with your own dataset [19].

For a future update of this script or development of a new one with the same topic you can explore this features.

## 8 Appendix

### 8.1 Exercise A - Solution

Solution for Exercise A, subsection 5.1.

Listing 65: Exercise A Solution

---

```
1  from openai import OpenAI
2  import pandas as pd
3  import re
4  import tiktoken
5  import PyPDF2
6
7  import ast
8  from scipy import spatial
9
10 client = openai.OpenAI()
11
12 SECTIONS_TO_IGNORE = [
13     "Contents",
14     "List of Tables",
15     "List of Figures",
16     "References",
17 ]
18
19 MAX_TOKENS = 1600
20 BATCH_SIZE = 1000
21
22 # TODO #1: Create consts for GPT_MODEL and EMBEDDING_MODEL (small)
23 EMBEDDING_MODEL = "text-embedding-3-small"
24 GPT_MODEL = "gpt-3.5-turbo"
25
26 def extract_text_from_pdf(pdf_path):
27     with open(pdf_path, 'rb') as file:
28         reader = PyPDF2.PdfReader(file)
29         text = ''
30         for page_num in range(len(reader.pages)):
31             text += reader.pages[page_num].extract_text()
32     return text
33
34 def split_sections_from_pdf(pdf_text):
35     title = []
36     text = []
37     ignore = True
38     current_section = "I.N.I.T.I.A.L-V.A.L.U.E"
39     for line in pdf_text.split('\n'):
40         line = line.strip()
41
42         if not line:
43             continue
44         if is_new_section(line):
45             ignore = ignore_section(line)
46             if ignore:
47                 continue
48             title.append(line)
49             if current_section != "I.N.I.T.I.A.L-V.A.L.U.E":
50                 text.append(current_section)
51             current_section = ""
52         else:
53             if not ignore:
54                 current_section += " " + line
```

```

55     if current_section:
56         text.append(current_section)
57
58
59     sections = [(title),(text)]
60     return sections
61
62 def ignore_section(line):
63     if any(section in line for section in SECTIONS_TO_IGNORE):
64         return True
65     return False
66
67 def is_new_section(line):
68     pattern = r"\d+\.\d+(?:\.\d+)? [A-Z].*?"
69
70     if line.strip().count(',') > 7:
71         return False
72     if re.match(pattern, line.strip()):
73         return True
74     if any(section in line for section in SECTIONS_TO_IGNORE):
75         return True
76
77     return False
78
79 def clean_section(section):
80
81     titles = section[0]
82     text = section[1]
83
84     for line in text:
85         line = re.sub(r"\[\d+\]", "", line)
86         line = re.sub(r"\[\d\d+\]", "", line)
87         line = line.strip()
88
89     return (titles, text)
90
91
92 def num_tokens(text, model = GPT_MODEL):
93     encoding = tiktoken.encoding_for_model(model)
94     return len(encoding.encode(text))
95
96 def halved_by_delimiter(string, delimiter = "\n"):
97     chunks = string.split(delimiter)
98     if len(chunks) == 1:
99         return [string, ""]
100    elif len(chunks) == 2:
101        return chunks
102    else:
103        total_tokens = num_tokens(string)
104        halfway = total_tokens // 2
105        best_diff = halfway
106        for i, chunk in enumerate(chunks):
107            left = delimiter.join(chunks[: i + 1])
108            left_tokens = num_tokens(left)
109            diff = abs(halfway - left_tokens)
110            if diff >= best_diff:
111                break
112            else:
113                best_diff = diff
114        left = delimiter.join(chunks[:i])
115        right = delimiter.join(chunks[i:])
116        return [left, right]
117

```

```

118
119 def truncated_string(string, model, max_tokens, print_warning = True,):
120     encoding = tiktoken.encoding_for_model(model)
121     encoded_string = encoding.encode(string)
122     truncated_string = encoding.decode(encoded_string[:max_tokens])
123     if print_warning and len(encoded_string) > max_tokens:
124         print(f"Warning: Truncated string from {len(encoded_string)} tokens to
125             {max_tokens} tokens.")
126
127
128 def split_strings_from_subsection(title, text, max_tokens = 1000, model = GPT_MODEL,
129     max_recursion = 5):
130
131     string = "\n\n".join(title + text)
132     num_tokens_in_string = num_tokens(string, model)
133
134     if num_tokens_in_string <= max_tokens:
135         return [string]
136
137     elif max_recursion == 0:
138         return [truncated_string(string, model, max_tokens)]
139
140     else:
141         for delimiter in ["\n\n", "\n", ". "]:
142             left, right = halved_by_delimiter(text, delimiter=delimiter)
143             if left == "" or right == "":
144                 continue
145             else:
146
147                 results = []
148                 for half in [left, right]:
149                     half_strings = split_strings_from_subsection(title,
150                         half,max_tokens,model,max_recursion - 1)
151                     results.extend(half_strings)
152
153         return results
154
155
156 # TODO #2: Call the previous created functions
157 pdf_text = extract_text_from_pdf("LETI_SISTCA_2023_24_Team2_OpenAI.pdf")
158
159 pdf_sections = split_sections_from_pdf(pdf_text)
160
161 cleaned_sections = clean_section(pdf_sections)
162
163 MAX_TOKENS = 1600
164 strings = []
165 titles = cleaned_sections[0]
166 texts = cleaned_sections[1]
167 for i in range(len(titles)):
168     strings.extend(split_strings_from_subsection(titles[i], texts[i],
169         max_tokens=MAX_TOKENS))
170
171
172 embeddings = []
173 for batch_start in range(0, len(strings), BATCH_SIZE):
174     batch_end = batch_start + BATCH_SIZE
175     batch = strings[batch_start:batch_end]
176     # TODO #3: Make a request to the embeddings API with the batch as input

```

```

177     response = client.embeddings.create(model=EMBEDDING_MODEL, input=batch)
178     for i, be in enumerate(response.data):
179         assert i == be.index
180     batch_embeddings = [e.embedding for e in response.data]
181     embeddings.extend(batch_embeddings)
182
183
184 df = pd.DataFrame({"text": strings, "embedding": embeddings})
185
186
187 SAVE_PATH = "SISTCA_TEAM2.csv"
188 df.to_csv(SAVE_PATH, index=False)
189
190 # TODO #4: Change the EMBEDDING_MODEL (ada)
191 EMBEDDING_MODEL = "text-embedding-ada-002"
192
193 GPT_MODEL = "gpt-3.5-turbo"
194
195 # TODO #5: Create a variable with the CSV file path
196 embeddings_path = "SISTCA_TEAM2.csv"
197
198 df = pd.read_csv(embeddings_path)
199 df['embedding'] = df['embedding'].apply(ast.literal_eval)
200
201 def strings_ranked_by_relatedness(query, df, relatedness_fn = lambda x, y: 1 -
202     spatial.distance.cosine(x, y), top_n = 100) :
203
204     # TODO 6: Make a request to the embeddings API with the query as input
205     query_embedding_response = client.embeddings.create(
206         model=EMBEDDING_MODEL,
207         input=query,
208     )
209
210     query_embedding = query_embedding_response.data[0].embedding
211     strings_and_relatednesses = [
212         (row["text"], relatedness_fn(query_embedding, row["embedding"]))
213         for i, row in df.iterrows()
214     ]
215     strings_and_relatednesses.sort(key=lambda x: x[1], reverse=True)
216     strings, relatednesses = zip(*strings_and_relatednesses)
217     return strings[:top_n], relatednesses[:top_n]
218
219     strings, relatednesses = strings_ranked_by_relatedness("open ai", df, top_n=5)
220
221 def num_tokens(text, model = GPT_MODEL):
222     encoding = tiktoken.encoding_for_model(model)
223     return len(encoding.encode(text))
224
225 def query_message(query,df, model, token_budget):
226     strings, relatednesses = strings_ranked_by_relatedness(query, df)
227     introduction = 'Use the below articles on the document about OpenAI made by Team 2,'
228     composed by Patricia Sousa, Carlos Alves, Jose Leal and Tiago Ribeiro, for SISTCA
229     to answer the subsequent question. If the answer cannot be found in the articles,
230     write "Sorry, the information you seek cannot be found in the document in
231     question."
232     question = f"\n\nQuestion: {query}"
233     message = introduction
234     for string in strings:
235         next_article = f'\n\nPDF article section:\n"""\n{string}\n"""'
236         if (
237             num_tokens(message + next_article + question, model)
238             > token_budget
239         ):

```

```

235         break
236     else:
237         message += next_article
238     return message + question
239
240
241 def ask(query, df = df, model = GPT_MODEL, token_budget = 4096 - 500, print_message =
242     False):
243     message = query_message(query, df, model=model, token_budget=token_budget)
244     if print_message:
245         print(message)
246     messages = [
247         {"role": "system", "content": "You answer questions about the document made by
248             Team2 for SISTCA about OpenAI."},
249         {"role": "user", "content": message},
250     ]
251
252     # TODO 7: Make a request to the Chat Completions API
253
254     response = client.chat.completions.create(
255         model=model,
256         messages=messages,
257         temperature=0
258     )
259     response_message = response.choices[0].message.content
260     return response_message
261
262 print(ask("Tell me about whisper or tts"))
263 print(ask("Give me the authors"))
264 print(ask("What can you tell me about the document"))
265 print(ask("Give me the document structure"))

```

---

## 8.2 Exercise B - Solution

Solution for Exercise B, subsection 5.2.

Listing 66: Exercise B Solution

```

1  from pathlib import Path
2  from openai import OpenAI
3
4  client = OpenAI()
5
6
7  # load the audio file
8  audio_input = open("Exercises/Exercise_B/audio.wav", "rb")
9
10
11 # transcribe contents and detect the input language
12 transcription = client.audio.transcriptions.create(
13     model="whisper-1",
14     file=audio_input,
15     response_format="verbose_json"
16 )
17
18 detected_language = transcription.language
19
20
21 # translate the transcription to another language using chatCompletions
22 output_language = "french"
23

```

```
24     translation = client.chat.completions.create(
25         model="gpt-3.5-turbo",
26         messages=[
27             {"role": "user", "content": f"Translate the following text to {output_language}:
28             {transcription.text}"}
29         ]
30     )
31 
32     translated_text = translation.choices[0].message.content
33 
34     # output the translated audio file
35     translated_audio = client.audio.speech.create(
36         model="tts-1",
37         voice="shimmer",
38         input=translated_text
39     )
40 
41     output_file_path = Path(__file__).parent / f"translation.mp3"
42 
43     translated_audio.write_to_file(output_file_path)
44 
45 
46     # test the output
47 
48     print(f"Detected language: {detected_language} \n")
49     print(f"Transcribed text: {transcription.text} \n")
50 
51     print(f"Translated text: {translated_text} \n")
52     print(f"Translated audio file saved to {output_file_path}")
```

---

## References

- [1] “Project Jupyter Documentation — Jupyter Documentation 4.1.1 alpha documentation.” [Online; accessed 2024-05-18].
- [2] IBM, “What is an ai model?,” *IBM*. Accessed: 2024-03-28.
- [3] R. Merritt, “What is a transformer model?,” *NVIDIA Blog*, mar 2022. (Accessed: 2024-03-28).
- [4] “Openai.” Available at: <https://openai.com/>. Accessed: 2024-03-27.
- [5] xAI, “xai grok.” Available at: <https://grok.x.ai/>. (Accessed: 2024-03-27).
- [6] Anthropic, “Anthropic home.” Available at: <https://www.anthropic.com/>. (Accessed: 2024-03-27).
- [7] Google, “Deepmind technologies.” Available at: <https://deepmind.google/technologies/>. (Accessed: 2024-03-27).
- [8] Cohere, “Cohere home.” Available at: <https://cohere.com/>. (Accessed: 2024-03-27).
- [9] “Python.” Available at: <https://wiki.python.org/>. (Accessed: 2024-03-29).
- [10] “Jupyter.” Available at: <https://jupyter.org/>. (Accessed: 2024-05-18).
- [11] “Hello gpt-4o.” Available at: <https://openai.com/index/hello-gpt-4o/>. (Accessed: 2024-05-18).
- [12] “Sora.” Available at: <https://openai.com/index/sora/>. (Accessed: 2024-05-18).
- [13] “Clip: Connecting text and images.” Available at: <https://openai.com/index/clip/>. (Accessed: 2024-05-18).
- [14] “Clip on github.” Available at: <https://github.com/openai/CLIP>. (Accessed: 2024-05-18).
- [15] “Jukebox sample explorer.” Available at: <https://jukebox.openai.com/>. (Accessed: 2024-05-18).
- [16] “Jukebox on github.” Available at: <https://github.com/openai/jukebox>. (Accessed: 2024-05-18).
- [17] “Point-e: A system for generating 3d point clouds from complex prompts.” Available at: <https://openai.com/index/point-e/>. (Accessed: 2024-05-18).
- [18] “Point-e on github.” Available at: <https://github.com/openai/point-e>. (Accessed: 2024-05-18).
- [19] “Fine-tuning.” Available at: <https://platform.openai.com/docs/guides/fine-tuning>. [Online; accessed 2024-05-18].