

# Learning Objectives

---

## Learners will be able to...

- Understand the importance of monitoring and analyzing API usage
- Learn how to use monitoring tools to track API usage and performance metrics such as response time, error rates, and usage patterns
- Understand how to analyze data collected from monitoring tools to identify areas for optimization and improvement
- Learn how to implement strategies to improve API performance and ensure efficient usage of resources

info

## Make Sure You Know

You are familiar with Python.

## Limitations

This is a gentle introduction. So there is a little bit of Python programming. The information might not be the most up to date as OpenAI releases new features.

# ChatGPT API Usage

As more applications incorporate the ChatGPT API by OpenAI, it becomes increasingly important to monitor and analyze API usage and performance to ensure optimal functionality and user experience. We are going to cover various tools and strategies for monitoring and analyzing the ChatGPT API, including using monitoring tools, analyzing data, and troubleshooting common issues.

To begin, let's look at an example of how to make an API call to the ChatGPT API using Python. For this example, we will use the "completion" feature of the ChatGPT API, which allows us to generate text based on a given prompt.

```
response=openai.ChatCompletion.create(
    model="gpt-3.5-turbo",
    messages=[
        {"role": "system", "content": "You are a helpful assistant."},
        {"role": "user", "content": "What's a good book to read on a rainy day? please just name 2."}
    ],
    n=3
)

for i in (response["choices"]):
    print(i["message"]['content'].strip())
    print(" ////////// ")
```

First of all we could use `openai.Completion.create` instead of `openai.ChatCompletion.create`. `openai.Completion.create` is used for generating text completions without a conversational context, while `openai.ChatCompletion.create` is used specifically for generating responses in a conversational context, making it easier to build chatbots and other conversational applications.

Let's quickly go over what the code does. The first import the necessary libraries (`os`, `openai`, and `secret`) and set the API key for the `openai` library using a separate file (`secret.api_key` contains the actual key).

We print the response from the ChatGPT API using a for loop to iterate over the choices list in the response object. Each choice object contains a text attribute, which contains the generated text from the ChatGPT API. We use the `strip()` method to remove any extra whitespace and add a separator (`/////////`) between each generated response.

# Time

Now that we've seen an example of how to use the openai library to make an API call to the ChatGPT API, let's explore some strategies for monitoring and analyzing API usage and performance to ensure optimal functionality and user experience.

One key metric to monitor is response **time**, which is the amount of time it takes for the API to generate a response to a given prompt. Long response times can lead to a poor user experience, so it's important to identify any bottlenecks in the API's processing pipeline and optimize them as needed. To monitor response time, you can use various tools such as profiling libraries like cProfile or built-in Python tools like the `time` module.

To monitor the response time of API requests to the ChatGPT API, you can use Python's built-in `time` module. Here's an example of how to measure the time it takes to generate a completion using the OpenAI library. First we are going to add the following library.

```
import time
```

Now we need to be selective with when we actually start timer. We are going to start the timer just before our API call.

```
# Start the timer
start_time = time.time()

# Make the API call to generate the completion
response=openai.ChatCompletion.create(
    model="gpt-3.5-turbo",
    messages=[
        {"role": "system", "content": "You are a helpful assistant."},
        {"role": "user", "content": "What's a good book to read on a rainy day? please just name 2."}
    ],
    n=3
)

# Calculate the elapsed time
elapsed_time = time.time() - start_time
print(f"Elapsed time: {elapsed_time:.4f} seconds")
```

We use the time module to measure the elapsed time between making the API call and receiving the response. We start the timer with `start_time = time.time()` and calculate the elapsed time with `elapsed_time = time.time() - start_time`. We then print the generated completion and elapsed time.

---

Another important metric to monitor is error rates, which are the number of errors that occur during API calls, such as timeouts or connection errors. High error rates can indicate issues with the API's infrastructure or other factors that are impacting performance. To monitor error rates, you can use log analysis tools like Loggly or Splunk, which can help you identify patterns in error messages and troubleshoot issues.

These tools allow you to analyze logs generated by your application and the ChatGPT API, which can help you identify patterns in error messages and troubleshoot issues.

To set up error logging with Loggly, you'll need to create an account and obtain an API key. You can then use the `loggly-python-handler` library to send log messages to Loggly.

```
from flask import Flask, request
import logging
import loggly.handlers

app = Flask(__name__)

# Set up the Loggly logger
log_handler = loggly.handlers.HTTPSHandler(
    'https://logs-01.loggly.com/inputs/YOUR-TOKEN-
    HERE/tag/python',
    'POST'
)
logger = logging.getLogger(__name__)
logger.addHandler(log_handler)

@app.route('/chat')
def chat():
    try:
        # ... code to handle API requests ...
    except Exception as e:
        logger.exception(f"API error: {e}")
        return "Error: " + str(e), 500
```

In this example, we set up a Flask application that handles API requests to the ChatGPT API. We use the `loggly` library to set up a logger that sends error messages to a Loggly endpoint. We set up the logger to use the

HTTPSHandler with the URL of the Loggly endpoint and the POST method to send the logs.

When an error occurs during an API request, we catch the exception using a try...except block and log the error message using `logger.exception()`. We then return an error message and a 500 status code to indicate that an error occurred.

# Optimize

## Optimizing ChatGPT API Performance

In addition to monitoring API usage and performance, it's also important to optimize the performance of the ChatGPT API itself. Here are some strategies for optimizing ChatGPT API performance:

- **Choose the Right Model**

The ChatGPT API offers a variety of models with different performance characteristics and capabilities. When building applications that use the ChatGPT API, it's important to choose the right model based on your specific needs and requirements. For example, if you need to generate responses quickly and with low latency, you may want to choose a smaller model with fewer parameters. On the other hand, if you need to generate more complex responses or handle a wider range of input types, you may want to choose a larger model with more parameters.

For that lets compare using gpt-4 and gpt-3.5 turbo with the prompt below. Run the code below twice, once with 3.5 and once with 4. . However GPT-4, excels at task requiring reason in comparison to gpt3.5 however, it will take longer to generates response. Please note at the time this course is being written the gpt-4 is not available yet for API calls

```
# Start the timer
start_time = time.time()

# Make the API call to generate the completion
response=openai.ChatCompletion.create(
    model="gpt-4-turbo",
    messages=[
        {"role": "system", "content": "You are a helpful assistant."},
        {"role": "user", "content": "What's a good book to read on a rainy day? please just name 2."}
    ],
    n=3
)

# Calculate the elapsed time
elapsed_time = time.time() - start_time
print(f"Elapsed time: {elapsed_time:.4f} seconds")
```

- **Input length**

The length of the input prompt can also impact the performance of the ChatGPT API. Longer prompts may require more processing time and resources, which can lead to longer response times and increased latency. To optimize performance, it's important to limit the length of input prompts and use concise, focused prompts that are optimized for the specific task or application.

- **Caching and preprocessing**

Caching and preprocessing are two techniques that can be used to optimize performance when using the ChatGPT API. Caching involves storing frequently used data in memory or on disk, so that it can be quickly retrieved without the need for additional processing.

Preprocessing involves analyzing input data and transforming it into a format that is optimized for the specific task or application, which can help reduce processing time and improve performance.

- **Scaling horizontally**

When building applications that use the ChatGPT API, it's important to consider scalability. As traffic to your application increases, you may need to scale horizontally by adding additional servers or resources to handle the increased load. By designing your application with scalability in mind and using cloud-based resources like AWS or Google Cloud, you can ensure that your application can handle increased traffic and continue to perform optimally.

# Securing and Protecting Data when using the ChatGPT API

When building applications that use the ChatGPT API, it's important to take steps to secure and protect the data being transmitted and processed. Throughout this course you can notice that we always interacted with our `api_key` in a secretive way that never exposes our `api_key` to the user.

```
import os
import openai
import secret
openai.api_key=secret.api_key
```

One of the first steps in securing the ChatGPT API is to ensure that API keys and other credentials are properly secured. API keys are used to authenticate API requests and must be kept secret to prevent unauthorized access. To secure API keys and credentials, it's important to use best practices such as encrypting secrets, storing secrets in secure locations, and limiting access to only authorized personnel. Here's an example of how to use Python's `dotenv` library to securely store API keys and other secrets:

```
import os
from dotenv import load_dotenv

# Load the environment variables from a .env file
load_dotenv()

# Retrieve the OpenAI API key from the environment variables
api_key = os.getenv("OPENAI_KEY")
```

We use the `dotenv` library to load environment variables from a `.env` file, which can be securely stored in a local directory. We then retrieve the OpenAI API key from the environment variables using `os.getenv()`, which prevents the API key from being exposed in plain text in the code.

---

Depending on your next project using the chatGPT API, it might be useful to add more layer of security.

Encrypt Data in Transit and at Rest:



To prevent data from being intercepted or stolen during transmission, it's important to use encryption when sending and receiving data to and from the ChatGPT API. Encryption can help protect data from being intercepted or stolen by hackers or other malicious actors. To encrypt data in transit, you can use HTTPS to encrypt traffic over the network.

To encrypt data at rest, you can use encryption libraries such as cryptography or pycryptodome to encrypt data before storing it in a database or file. Here's an example of how to use the cryptography library to encrypt data at rest:

```
from cryptography.fernet import Fernet

# Generate a secret key for encryption
key = Fernet.generate_key()

# Create a Fernet object with the secret key
fernet = Fernet(key)

# Encrypt data using the Fernet object
encrypted_data = fernet.encrypt(b"Hello, world!")

# Decrypt data using the Fernet object
decrypted_data = fernet.decrypt(encrypted_data)
```

In this example, we use the cryptography library to generate a secret key for encryption using `Fernet.generate_key()`. We then create a Fernet object with the secret key, which can be used to encrypt and decrypt data using the `encrypt()` and `decrypt()` methods. By encrypting data at rest, you can help prevent unauthorized access to sensitive information.

# Coding Exercise

In this exercise, we'll be writing a function `Generate` which uses OpenAI's API to generate text. This function will take a list of prompts from the user, process each prompt with the OpenAI model, and return the response generated by each prompt, as well as the time it took to generate each response.

Here's how the `Generate` function works:

- The function accepts a list of prompts as its parameter.
- For each prompt, we record the start time and then use `openai.Completion.create` to generate a response from OpenAI. We limit the response to 100 tokens, but you can adjust this limit based on your needs.
- We then record the end time and calculate the time taken for the generation.
- We store the response text and the time taken in a tuple, and append this tuple to our list of responses.
- Finally, we return the list of responses.

You can call this function with a list of prompts like this:

```
prompts = ["Who won the world series in 2020?", "What's the capital of France?", "Who wrote Pride and Prejudice?"]
responses = Generate(prompts)

for i, (response, time_taken) in enumerate(responses):
    print(f"Prompt: {prompts[i]}\nResponse: {response}\nTime Taken: {time_taken} seconds\n")
```

The response generated should look like the following:

```
Prompt: Who won the world series in 2020?
Response: The Los Angeles Dodgers won the 2020 World Series.
Time Taken: 0.42620301246643066 seconds

Prompt: What's the capital of France?
Response: Paris
Time Taken: 0.3195617198944092 seconds

Prompt: Who wrote Pride and Prejudice?
Response: Jane Austen
Time Taken: 0.39693474769592285 seconds
```

