

GENAI - CONTAINERIZED VIDEO TRANSCRIPTION AND CHAT

Week 11 Homework

CourseCloud Computing Infrastructure

Name: Hartina Vonyee Cleon

ID: 20145

Due Date: August 3, 2024

Introduction

In this task, I deployed the Streamlit application inside a Docker container on Google Cloud Kubernetes. The aim was to set up the necessary environment, obtain and apply various API keys, create the Docker image, and ensure the application functions optimally within a containerized environment. I also addressed how to handle issues such as Docker Compose file validation, environment variables, and obtaining information about specific containers, such as IP addresses. By following a planned approach to each design step, I aimed to create a functional, cloud-based application. Specifically, this project involves designing a chatbot that can respond to questions from a video. It demonstrates the integration of technologies such as Docker, OpenAI, Whisper, Embeddings, Chat completions, Pinecone, and Retrieval-Augmented Generation.

Prerequisites

- OpenAI API Key
- Pinecone API Key
- The latest version of Docker Desktop
- GitHub repository

Steps to Build and Run the Application

1. Clone Repository: Use Git to clone the project's repository.

Now that we have all the necessary resources for the task, we can proceed with developing our application. In the terminal, navigate to the desired directory:

```
git clone https://github.com/Davidnet/docker-genai.git
```

```
C:\Users\HP\Documents\SEM 2\CS571>clone https://github.com/Davidnet/docker-genai.git
'clone' is not recognized as an internal or external command,
operable program or batch file.

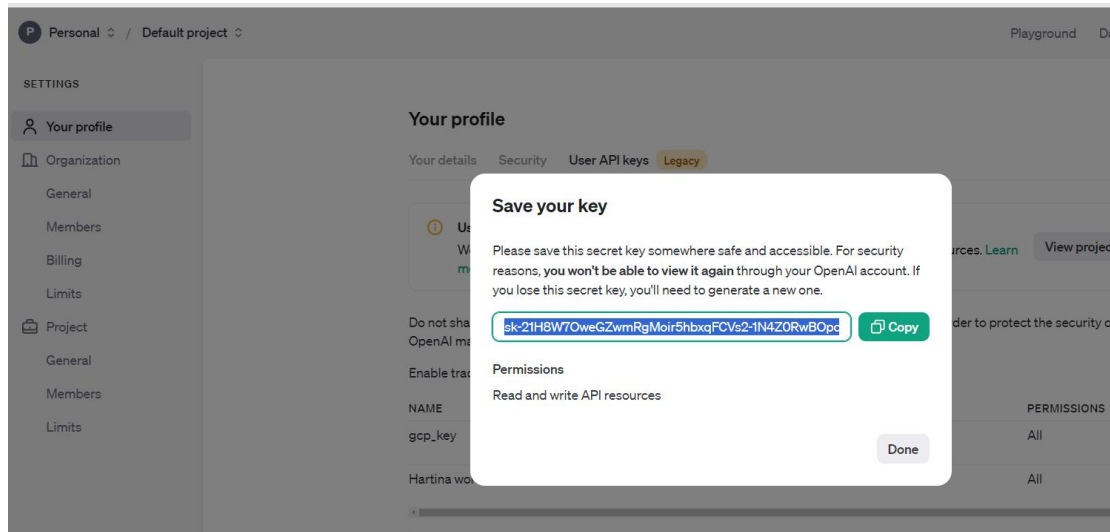
C:\Users\HP\Documents\SEM 2\CS571>git clone https://github.com/Davidnet/docker-genai.git
Cloning into 'docker-genai'...
remote: Enumerating objects: 66, done.
remote: Counting objects: 100% (66/66), done.
remote: Compressing objects: 100% (43/43), done.
remote: Total 66 (delta 24), reused 60 (delta 20), pack-reused 0
Receiving objects: 100% (66/66), 114.38 KiB | 329.00 KiB/s, done.
Resolving deltas: 100% (24/24), done.
```

Get an OpenAI API Key

To obtain an OpenAI API key, follow these steps:

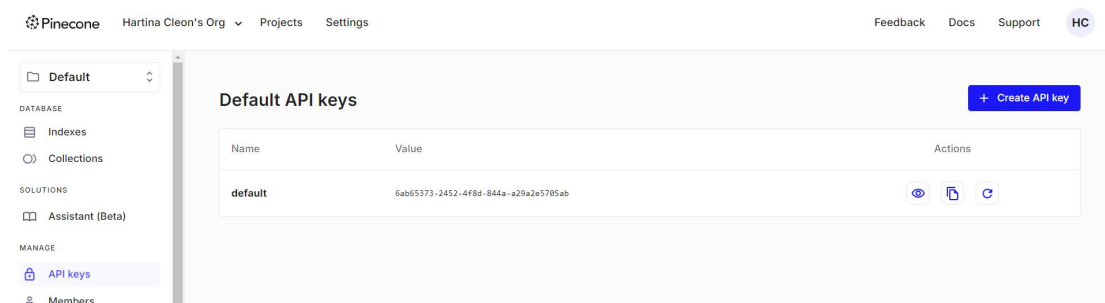
1. Click on the provided link to access the OpenAI website.
2. Enter your OpenAI login details.
3. Navigate to the 'Billing' section.
4. Click on 'Add Payment Details' and proceed to add the required credit.
5. Once you have successfully added enough credit, go to the sidebar menu.
6. From the menu, select "Create a New Secret Key."

7. Optionally, you can provide a name for this key.
8. Fill in the necessary information in the "Create Secret Key" section to generate your API key.



Obtain a Pinecone API Key

To get a Pinecone API Key, follow the link provided. This will take you to the Pinecone web app where you can either enter your account credentials or create a new account. The link will open in a new window. Go to the side menu and navigate to API KEYS to retrieve your default API Key.



Specify your API keys

Go to docker-genai directory:

```
cd docker-genai
```

Create a text file ".env":

```
code .env
```

Give the below script in the .env file and replace the keys with your Actual keys:

Specify your API keys in ".env" using this template:

```
#-----  
# OpenAI
```

```
#-----
OPENAI_TOKEN=your-api-key # Replace your-api-key with your personal API key
#-----
# Pinecone
#-----
PINECONE_TOKEN=your-api-key # Replace your-api-key with your personal API key
```

Build and run the application

docker compose up - --build

You will notice this while Docker is constructing the application; the following is typical when the application is up. As you can deduce, the Dockerfile is similar to the one used when running the application.

```
C:\Users\HP\Documents\SEM 2\CS571\docker-genai>docker compose up --build
2024/08/04 02:20:41 http2: server: error reading preface from client //./pipe/docker_engine: file has already been closed
[*] Building 0.0s (0/0) docker:default
[*] Building 0.0s (0/0) docker:defaulttr reading preface from client //./pipe/docker_engine: file has already been closed
[*] Building 63.0s (20/20) FINISHED
=> [yt-whisper internal] load build definition from Dockerfile
=> transferring dockerfile: 1.88kB
=> [bot internal] load build definition from Dockerfile
=> transferring dockerfile: 1.88kB
=> [bot] resolve image config for docker-image://docker.io/docker/dockerfile:1
=> [bot auth] docker/dockerfile:pull token for registry-1.docker.io
=> CACHED [yt-whisper internal] docker-image://docker.io/docker/dockerfile:1@sha256:fe40cf4e92cd0c467be2cfc30657a680ae2398318afd50b0c80585784c604f28
=> [bot internal] load metadata for docker.io/library/python:3.11-slim
=> [bot auth] library/python:pull token for registry-1.docker.io
=> [bot internal] load .dockerignore
=> transferring context: 2B
=> [yt-whisper internal] load .dockerignore
=> transferring context: 2B
=> [yt-whisper base 1/5] FROM docker.io/library/python:3.11-slim@sha256:7f49f147e57a65a5ca731203ed358ac5c88fa54aeb942924dd7057fe34a18e79
=> resolve docker.io/library/python:3.11-slim@sha256:7f49f147e57a65a5ca731203ed358ac5c88fa54aeb942924dd7057fe34a18e79
=> sha256:446018211442caf50e259f716bbf6a79c1712186fde9fa83772fa081d6ada8dc 3.51MB / 3.51MB
=> sha256:5bcb0eaa0b9f2d9aa541cc89c71c5fb6fc2e5a981ce3e9c9785aba1311aa9551 12.87MB / 12.87MB
=> sha256:1223fada08950a6ff83d6d4f82381992c60b6b18a3e46677968b8a4e29408f2b 233B / 233B
=> sha256:7f49f147e57a65a5ca731203ed358ac5c88fa54aeb942924dd7057fe34a18e79 9.12kB / 9.12kB
=> sha256:7adf7ab12f5e74ac929e77c289017c8918026d2127b9a7b8da8a8447986732d3 1.94kB / 1.94kB
=> sha256:58a64f49c6712896f9ee05afe5ac2c39568951dc82ab11b9a8c709755b9b55a1 6.92kB / 6.92kB
=> sha256:ealc5d8bd5c8b355169216f80bdac84e7c0f954b94dc6570d0ca23c448c414d3 3.21MB / 3.21MB
=> extracting sha256:446018211442caf50e259f716bbf6a79c1712186fde9fa83772fa081d6ada8dc 0.25
=> extracting sha256:5bcb0eaa0b9f2d9aa541cc89c71c5fb6fc2e5a981ce3e9c977968b8a4e29408f2b 0.55
=> extracting sha256:1223fada08950a6ff83d6d4f82381992c60b6b18a3e46677968b8a4e29408f2b 0.05
=> extracting sha256:ealc5d8bd5c8b355169216f80bdac84e7c0f954b94dc6570d0ca23c448c414d3 0.35
=> [bot internal] load build context
=> transferring context: 132.54kB
=> [yt-whisper internal] load build context
=> exporting layers
=> writing image sha256:59836ce5c8007517fe16a9d8ce87d8820cd915732de13c6acfcald8a899c980e2
=> naming to docker.io/library/docker-genai-bot
[+] Running 3/3
✔ Network docker-genai_default Created
✔ Container docker-genai-yt-whisper-1 Created
✔ Container docker-genai-bot-1 Created
Attaching to bot-1, yt-whisper-1
yt-whisper-1 |
yt-whisper-1 | Collecting usage statistics. To deactivate, set browser.gatherUsageStats to false.
yt-whisper-1 |
bot-1 |
bot-1 | Collecting usage statistics. To deactivate, set browser.gatherUsageStats to false.
yt-whisper-1 |
bot-1 |
yt-whisper-1 | You can now view your Streamlit app in your browser.
bot-1 |
yt-whisper-1 | You can now view your Streamlit app in your browser.
yt-whisper-1 | URL: http://0.0.0.0:8503
bot-1 |
yt-whisper-1 |
bot-1 | URL: http://0.0.0.0:8504
bot-1 |
```

Use the yt-whisper service

Open a browser and access the yt-whisper service at <http://localhost:8503>.
Enter the Youtube video URL you want to use and select “Submit”:


```
    return self._request(
        ^^^^^^^^^^^^^^^^^^^
File "/usr/local/lib/python3.11/site-packages/openai/_base_client.py", line 1031, in _request
    return self._retry_request(
        ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
File "/usr/local/lib/python3.11/site-packages/openai/_base_client.py", line 1079, in _retry_request
    return self._request(
        ^^^^^^^^^^^^^^^^^^^
File "/usr/local/lib/python3.11/site-packages/openai/_base_client.py", line 1031, in _request
    return self._retry_request(
        ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
File "/usr/local/lib/python3.11/site-packages/openai/_base_client.py", line 1079, in _retry_request
    return self._request(
        ^^^^^^^^^^^^^^^^^^^
File "/usr/local/lib/python3.11/site-packages/openai/_base_client.py", line 1046, in _request
    raise self._make_status_error_from_response(err.response) from None
```