**Step 1: GenAI - Containerize your app.**

1. First, install the latest version of Docker Desktop for windows.

   - [Docker Desktop: The #1 Containerization Tool for Developers | Docker](#)



2. Go to the terminal and navigate to our working directory.

```
C:\Users\patel>mkdir GenAIApplication

C:\Users\patel>cd GenAIApplication

C:\Users\patel\GenAIApplication>
```

3. Clone the sample application. We run the following command to clone the repository:

git clone [https://github.com/craig-osterhout/docker-genai-sample](https://github.com/craig-osterhout/docker-genai-sample)

```
C:\Users\patel\GenAIApplication>git clone https://github.com/craig-osterhout/docker-genai-sample
Cloning into 'docker-genai-sample'...
remote: Enumerating objects: 11, done.
remote: Counting objects: 100% (11/11), done.
remote: Compressing objects: 100% (10/10), done.
remote: Total 11 (delta 0), reused 11 (delta 0), pack-reused 0 (from 0)
Receiving objects: 100% (11/11), 10.17 KiB | 5.08 MiB/s, done.

C:\Users\patel\GenAIApplication>cd docker-genai-sample
```

4. You should now have the following files in your docker-genai-sample directory.

```
C:\Users\patel\GenAIApplication\docker-genai-sample>dir
 Volume in drive C is OS
 Volume Serial Number is 8AA4-9ED7

 Directory of C:\Users\patel\GenAIApplication\docker-genai-sample

11/27/2024  06:24 PM    <DIR>          .
11/27/2024  06:24 PM    <DIR>          ..
11/27/2024  06:24 PM             3,895 app.py
11/27/2024  06:24 PM             9,099 chains.py
11/27/2024  06:24 PM               967 env.example
11/27/2024  06:24 PM             7,169 LICENSE
11/27/2024  06:24 PM               179 README.md
11/27/2024  06:24 PM               106 requirements.txt
11/27/2024  06:24 PM             1,945 utils.py
               7 File(s)         23,360 bytes
               2 Dir(s)  880,910,508,032 bytes free
```

5. Now that we have an application, we can use docker init to create the necessary Docker assets to containerize our application. Inside the docker-genai-sample directory, run the docker init command.

```
C:\Users\patel\GenAIApplication\docker-genai-sample>docker init

Welcome to the Docker Init CLI!

This utility will walk you through creating the following files with sensible defaults for your project:
  - .dockerignore
  - Dockerfile
  - compose.yaml
  - README.Docker.md

Let's get started!

? What application platform does your project use? Python
? What version of Python do you want to use? 3.11.7

? What version of Python do you want to use? 3.11.7
? What port do you want your app to listen on? (8000) 8000

? What port do you want your app to listen on? 8000
? What is the command you use to run your app (e.g., gunicorn 'myapp.example:app' --bind=0.0.0.0:8000)? streamlit run app.py --server.address=0.0.0 --server.port=8000
? What is the command you use to run your app (e.g., gunicorn 'myapp.example:app' --bind=0.0.0.0:8000)? streamlit run app.py --server.address=0.0.0 --server.port=8000

✔ Created → .dockerignore
✔ Created → Dockerfile
✔ Created → compose.yaml
✔ Created → README.Docker.md

→ Your Docker files are ready!
  Review your Docker files and tailor them to your application.
  Consult README.Docker.md for information about using the generated files.

What's next?
  Start your application by running → docker compose up --build
  Your application will be available at http://localhost:8000

C:\Users\patel\GenAIApplication\docker-genai-sample>
```

**Step 2: GenAI - Develop your app.**
Adding a Local Database
Here we will update the compose.yaml file to define a database service, and we will specify an environment variables file to load the database connection

information rather than manually entering the information every time. To run the database service:

1. In the cloned repository's directory, rename env.example file to .env. This file contains the environment variables that the containers will use.

```
C:\Users\patel\GenAIApplication\docker-genai-sample>ren env.example .env

C:\Users\patel\GenAIApplication\docker-genai-sample>dir
 Volume in drive C is OS
 Volume Serial Number is 8AA4-9ED7

 Directory of C:\Users\patel\GenAIApplication\docker-genai-sample

11/27/2024  07:06 PM    <DIR>          .
11/27/2024  06:24 PM    <DIR>          ..
11/27/2024  06:44 PM               629 .dockerignore
11/27/2024  06:24 PM               967 .env
11/27/2024  06:24 PM             3,895 app.py
11/27/2024  06:24 PM             9,099 chains.py
11/27/2024  06:44 PM             1,642 compose.yaml
11/27/2024  06:44 PM             1,667 Dockerfile
11/27/2024  06:24 PM             7,169 LICENSE
11/27/2024  06:44 PM               826 README.Docker.md
11/27/2024  06:24 PM               179 README.md
11/27/2024  06:24 PM               106 requirements.txt
11/27/2024  06:24 PM             1,945 utils.py
              11 File(s)         28,124 bytes
               2 Dir(s)  880,880,824,320 bytes free

C:\Users\patel\GenAIApplication\docker-genai-sample>
```

2. Then open the compose.yaml file in an IDE or text editor.
   - Add instructions to run a Neo4j database.
   - Specify the environment file under the server service in order to pass in the environment variables for the connection.

```
services:
  server:
    build:
      context: .
    ports:
      - "8000:8000"
    env_file:
      - .env
```

```
    depends_on:
      database:
        condition: service_healthy

  database:
    image: neo4j:5.11
    ports:
      - "7474:7474"
      - "7687:7687"
    environment:
      - NEO4J_AUTH=${NEO4J_USERNAME}/${NEO4J_PASSWORD}
    healthcheck:
      test: ["CMD-SHELL", "wget --no-verbose --tries=1 --spider
localhost:7474 || exit 1"]
      interval: 5s
      timeout: 3s
      retries: 5
```

4. Run the application. Inside the docker-genai-sample directory, run the following command in a terminal.
Before running below step, open your Docker Desktop
 -   wsl --list --verbose

```
PS C:\Users\patel\GenAIApplication\docker-genai-sample> wsl --list --verbose
  NAME                STATE           VERSION
* docker-desktop      Running         2
```

 -   docker compose up --build

```
C:\Users\patel\GenAIApplication\docker-genai-sample>docker compose up --build
time="2024-11-27T19:22:16-08:00" level=warning msg="C:\\Users\\patel\\GenAIApplication\\docker-genai-sample\\c
ompose.yaml: the attribute `version` is obsolete, it will be ignored, please remove it to avoid potential conf
usion"
[+] Running 6/6
 ✔database Pulled                                                                           34.0s
   ✔33a66ada74dc Download complete                                                          27.1s
   ✔732d09690fed Download complete                                                          30.2s
   ✔e8cba66f5b65 Download complete                                                           0.6s
   ✔7d97e254a046 Download complete                                                           9.0s
   ✔9e41d761a8cf Download complete                                                           0.6s
[+] Building 144.5s (9/11)                                                     docker:desktop-linux
 => [server] resolve image config for docker-image://docker.io/docker/dockerfile:1            1.6s
 => [server] docker-image://docker.io/docker/dockerfile:1@sha256:865e5dd094beca432e8c0a1d5e1c465db5f998  1.7s
```

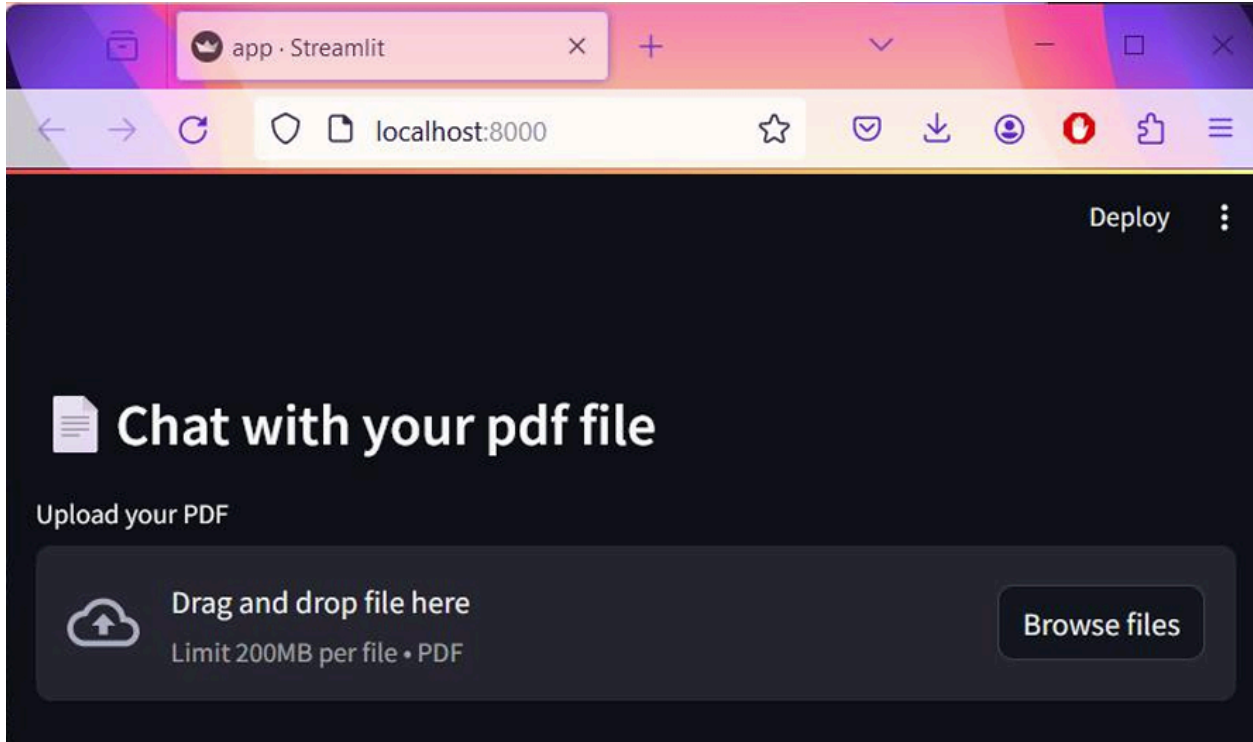 -   We can also see the progress from Docker Desktop.

5. Access the application. Open a browser and view the application at http://localhost:8000. You should see a simple Streamlit application.

6. Stop the application. In the terminal, press ctrl+c to stop the application.



**Adding a Local or Remote LLM Service**
**1. Install the prerequisites.**
- For Docker Engine on Linux, install the NVIDIA Container Toolkilt.
- For Docker Desktop on Windows 10/11, install the latest **NVIDIA driver** and make sure you are using the **WSL2 backend**

6

**2. Add the Ollama service and a volume in your compose.yaml. The following is the updated compose.yaml:**

```yaml
version: "3.8"

services:
  server:
    build:
      context: .
    ports:
      - "8000:8000"
    env_file:
      - .env
    depends_on:
      database:
        condition: service_healthy

  database:
    image: neo4j:5.11
    ports:
      - "7474:7474"
      - "7687:7687"
    environment:
      - NEO4J_AUTH=${NEO4J_USERNAME}/${NEO4J_PASSWORD}
    healthcheck:
      test: ["CMD-SHELL", "wget --no-verbose --tries=1 --spider localhost:7474 || exit 1"]
      interval: 5s
      timeout: 3s
      retries: 5
  ollama:
   image: ollama/ollama:latest
   ports:
     - "11434:11434"
   volumes:
     - ollama_volume:/root/.ollama
   deploy:
    resources:
      reservations:
        devices:
          - driver: nvidia
            count: all
            capabilities: [gpu]
  volumes:
    ollama_volume:
```

3. Add the ollama-pull service to your compose.yaml file. This service uses the docker/genai:ollama-pull image, based on the GenAI Stack's pull_model.Dockerfile and will automatically pull the model for your Ollama container. The following is the updated section of the compose.yaml file:

```yaml
version: "3.8"

services:
  server:
    build:
      context: .
    ports:
      - "8000:8000"
    env_file:
      - .env
    depends_on:
      database:
        condition: service_healthy
      ollama-pull:
        condition: service_completed_successfully
  ollama-pull:
    image: docker/genai:ollama-pull
    env_file:
      - .env

  database:
    image: neo4j:5.11
    ports:
      - "7474:7474"
      - "7687:7687"
    environment:
      - NEO4J_AUTH=${NEO4J_USERNAME}/${NEO4J_PASSWORD}
    healthcheck:
      test: ["CMD-SHELL", "wget --no-verbose --tries=1 --spider localhost:7474 || exit 1"]
      interval: 5s
      timeout: 3s
      retries: 5
  ollama:
   image: ollama/ollama:latest
   ports:
     - "11434:11434"
   volumes:
     - ollama_volume:/root/.ollama
   deploy:
    resources:
      reservations:
        devices:
          - driver: nvidia
            count: all
            capabilities: [gpu]
  volumes:
    ollama volume:
```

2. Update the OLLAMA_BASE_URL value in your .env file to
   http://host.docker.internal:11434

```
#*********************************************************************
# LLM and Embedding Model
#*********************************************************************
LLM=llama2 # Set to "gpt-3.5" to use OpenAI.
EMBEDDING_MODEL=sentence_transformer


#*********************************************************************
# Neo4j
#*********************************************************************
NEO4J_URI=neo4j://database:7687
NEO4J_USERNAME=neo4j
NEO4J_PASSWORD=password


#*********************************************************************
# Ollama
#*********************************************************************
OLLAMA_BASE_URL=http://host.docker.internal:11434        #http://ollama:11434


#*********************************************************************
# OpenAI
#*********************************************************************
# Only required when using OpenAI LLM or embedding model
# OpenAI charges may apply. For details, see
# https://openai.com/pricing

#OPENAI_API_KEY=sk-..
```
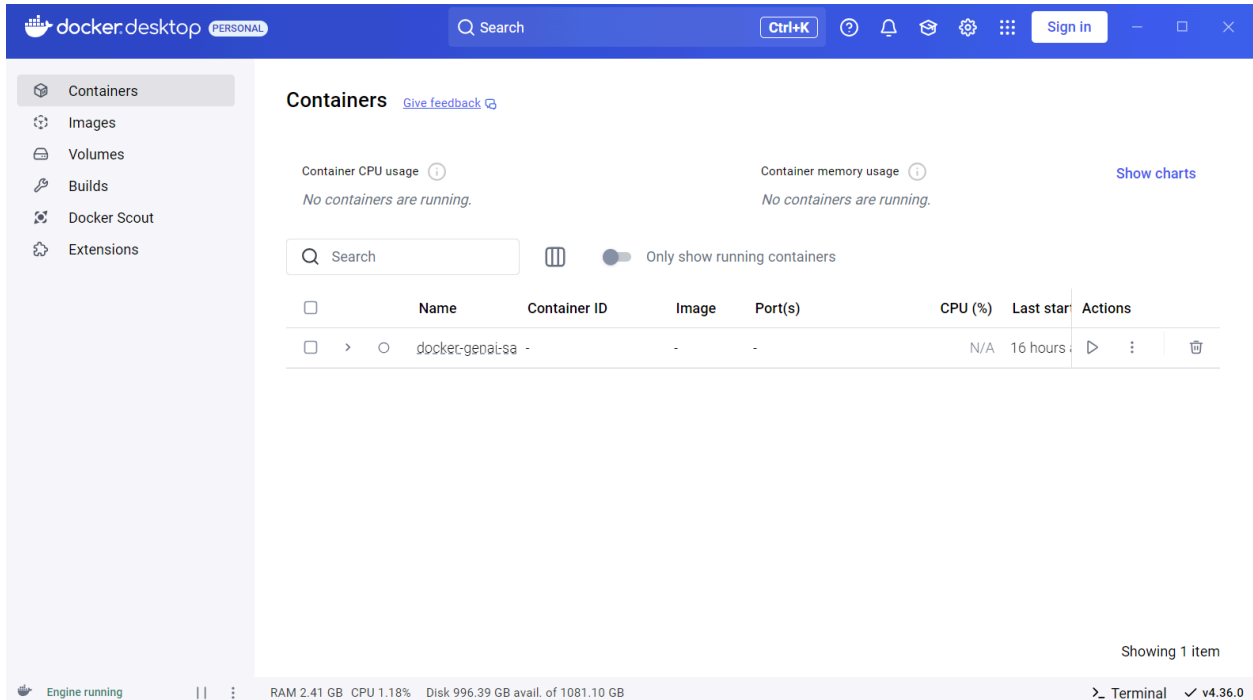
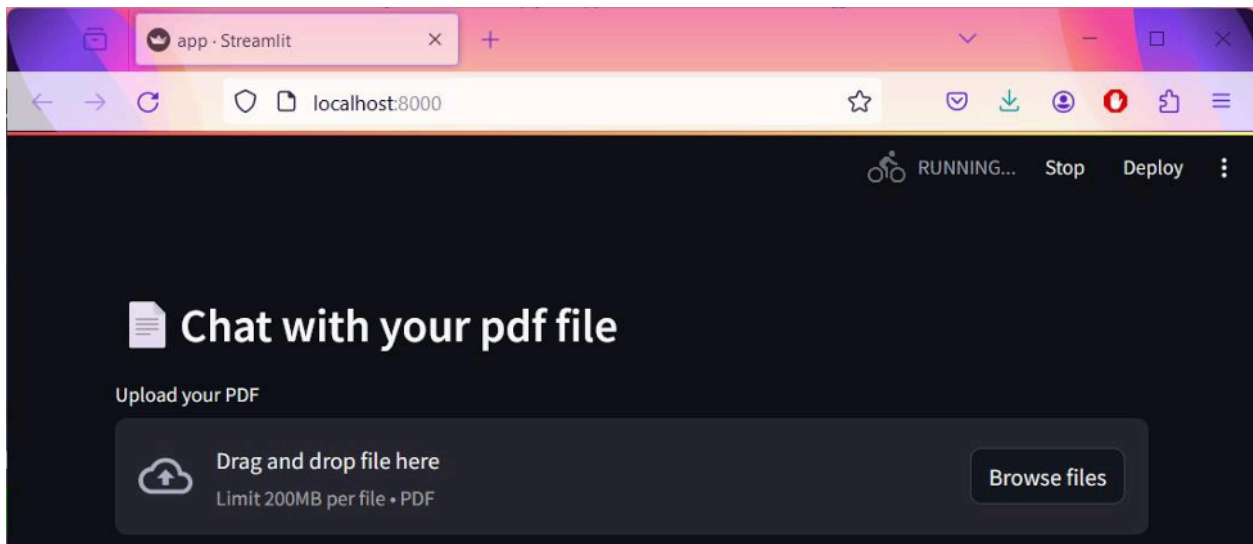## Add configuration for ollama pull and ollam in docker compose file

```
C:\Users\patel\GenAIApplication\docker-genai-sample>docker compose up --build
time="2024-11-28T11:23:11-08:00" level=warning msg="C:\\Users\\patel\\GenAIApplication\\docker-genai-sample\\compose.yaml: the attribute `versio
n` is obsolete, it will be ignored, please remove it to avoid potential confusion"
[+] Running 5/13
 - ollama-pull [▮▮▮▮▮] 21.69MB / 438.5MB Pulling                                                          54.9s
   ✔ d7f704120c50 Download complete                                                                        0.7s
   - a48641193673 Downloading [=========>                 ]  6.291MB/29.55MB                               52.7s
   ✔ 496e8c35aa41 Download complete                                                                       19.4s
   ✔ 4f4fb700ef54 Download complete                                                                        0.7s
   ✔ 1eadfce5a711 Download complete                                                                        0.8s
   - 25219de7956a Downloading [============>              ]  6.291MB/25.6MB                                52.7s
   - 0977b56ccc02 Downloading [>                          ]  6.291MB/380.6MB                               52.7s
 - ollama [▮▮] 22.02MB / 1.888GB Pulling                                                                  54.9s
   - b488f0047914 Downloading [>                          ]  9.437MB/1.848GB                               52.7s
   ✔ f46f4148708e Download complete                                                                        4.0s
   - 64l4378b6477 Downloading [=========>                 ]  6.291MB/29.54MB                               52.7s
   - 719d99f741d7 Downloading [==============================>]  6.291MB/9.694MB                           52.7s
```
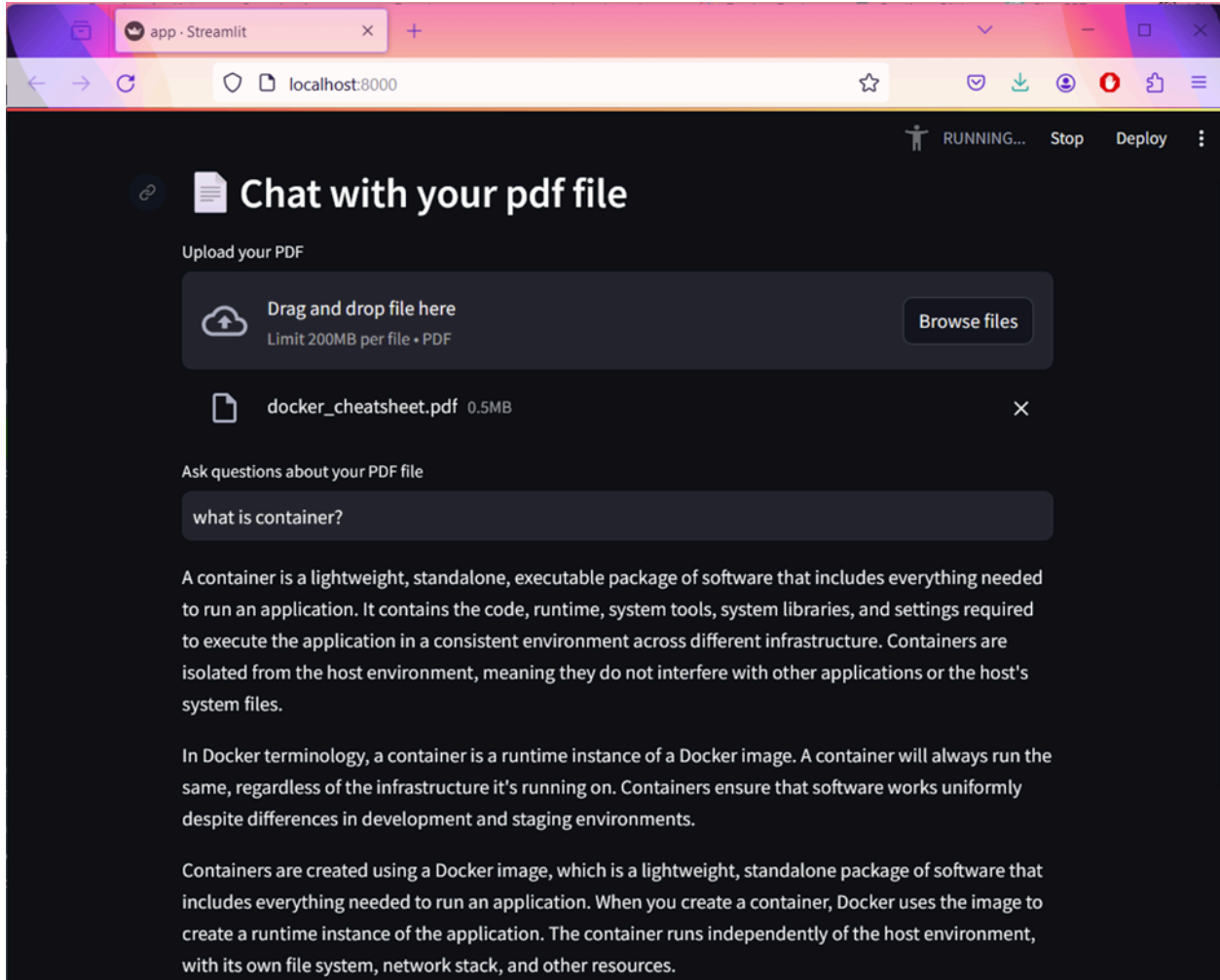
**Build new images and containers using docker compose --build**

**2. Once the application is running, open a browser and access the application at http://localhost:8000.**



**3. Then we can upload a PDF file, for example the Docker CLI Cheat Sheet, and ask a question about the PDF.**

Through this we have set up a development environment that provides access to all the services that our GenAI application needs.

**Link to GitHub - [Cloud-Computing/kubernetes at main · hpatel65373/Cloud-Computing](#)**