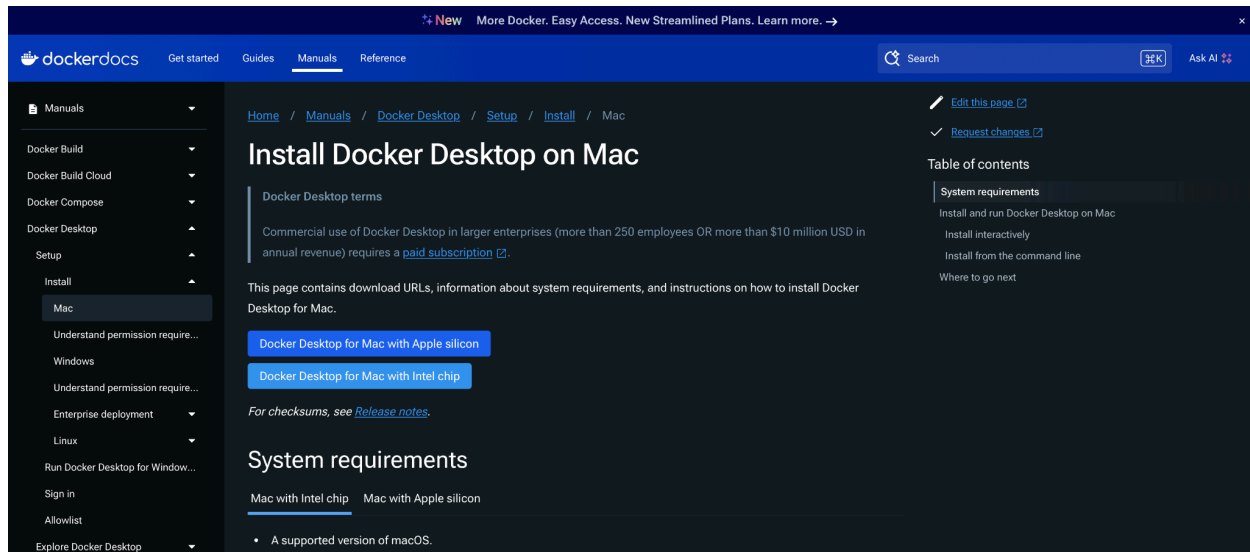Week 11 Homework 2: Project: GenAI - Develop your containerized app
## Step 1: GenAI - Containerize your app

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



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

3. Clone the sample application. We run the following command to clone the repository:
git clone https://github.com/craig-osterhout/docker-genai-sample

```
rashmipurandare@Rashmis-Laptop ~ % mkdir GenAIApp
rashmipurandare@Rashmis-Laptop ~ % cd GenAIApp
rashmipurandare@Rashmis-Laptop GenAIApp % 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 | 2.54 MiB/s, done.
rashmipurandare@Rashmis-Laptop GenAIApp %
```

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

```
rashmipurandare@Rashmis-Laptop GenAIApp % cd docker-genai-sample
rashmipurandare@Rashmis-Laptop docker-genai-sample % ls
LICENSE               app.py               env.example          utils.py
README.md             chains.py            requirements.txt
```

4. 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.

Week 11 Homework 2: Project: GenAI - Develop your containerized app

```
  — .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? 3.13.0
? 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=0.0.0.0 --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
```

## 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.

```
[rashmipurandare@Rashmis-Laptop docker-genai-sample % mv env.example .env
```

2. Then open the compose.yaml file in an IDE or text editor.
$ nano compose.yaml

3. In the compose.yaml file, add the following:
          o Add instructions to run a Neo4j database
          o 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:

## Week 11 Homework 2: Project: GenAI - Develop your containerized app

```
    - .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.
$ docker compose up --build

Week 11 Homework 2: Project: GenAI - Develop your containerized app

• 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.

• Note that asking questions to a PDF will cause the application to fail because the
LLM service specified in the .env file isn't running yet.



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

Week 11 Homework 2: Project: GenAI - Develop your containerized app
**Adding a Local or Remote LLM Service**

1. Install the prerequisites.
• For Docker Engine on Linux, install the NVIDIA Container Toolkit.
• For Docker Desktop on Windows 10/11, install the latest NVIDIA driver and make sure you are using the WSL2 backend

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

```
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:

```
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
```

**Run Ollama in a Container**

1. Install and run Ollama on your host machine.

# Welcome to Ollama

Let's get you up and running with
your own large language models.

**Next**

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

Week 11 Homework 2: Project: GenAI - Develop your containerized app

```
  UW PICO 5.09                                                    File: .env

  #*******************************************************************
  # 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

  #*******************************************************************
  # 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-..
```

3. Pull the model to Ollama using the following command.
$ ollama pull llama2

```
[rashmipurandare@Rashmis-Laptop docker-genai-sample % nano .env
[rashmipurandare@Rashmis-Laptop docker-genai-sample % ollama pull llama2
pulling manifest
pulling 8934d96d3f08... 100% ▐████████████████████████████▌  3.8 GB
pulling 8c17c2ebb0ea... 100% ▐████████████████████████████▌  7.0 KB
pulling 7c23fb36d801... 100% ▐████████████████████████████▌  4.8 KB
pulling 2e0493f67d0c... 100% ▐████████████████████████████▌   59 B
pulling fa304d675061... 100% ▐████████████████████████████▌   91 B
pulling 42ba7f8a01dd... 100% ▐████████████████████████████▌  557 B
verifying sha256 digest
writing manifest
success
```

Note: In case you are using OpenAI you can do the following steps instead.
i. ii. Update the LLM value in your .env file to gpt-3.5.
Uncomment and update the OPENAI_API_KEY value in your .env file to your OpenAI
API key


**Run Your GenAI Application**


1. To run all the services, run the following command.
$ docker compose up --build

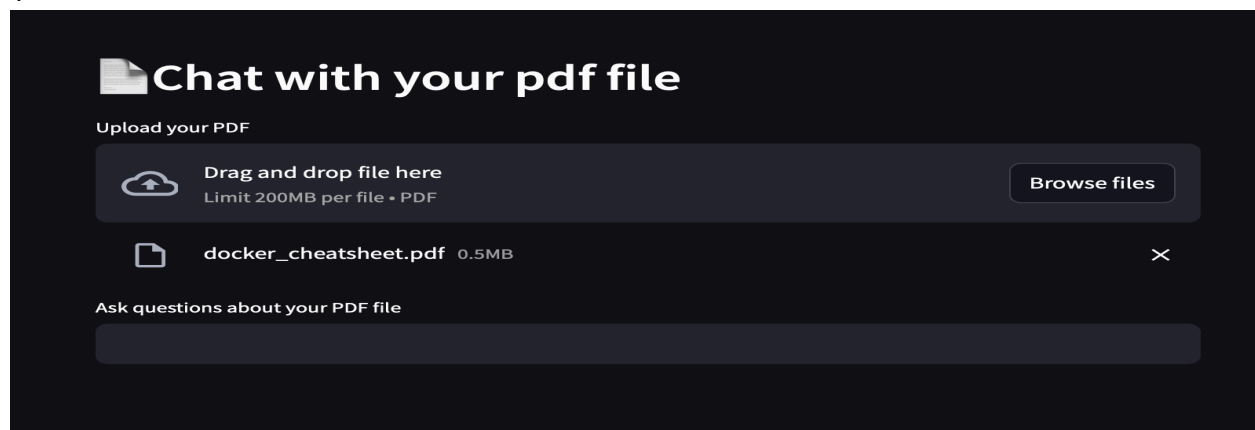Week 11 Homework 2: Project: GenAI - Develop your containerized app

```
=> [server internal] load build definition from Dockerfile
=> => transferring dockerfile: 1.59kB
=> [server] resolve image config for docker-image://docker.io/docker/dockerfile:1
=> CACHED [server] docker-image://docker.io/docker/dockerfile:1@sha256:865e5dd094beca432e8c0a1d5e1c465db5f998dca4e439981029b3b81fb3
=> [server internal] load metadata for docker.io/library/python:3.11-slim
=> [server internal] load .dockerignore
=> => transferring context: 708B
=> [server stage-0 1/5] FROM docker.io/library/python:3.11-slim@sha256:e8381c802593deb0c4d25bd3f4e05e94382f6bf33090de22679fc7488cde
=> [server internal] load build context
=> => transferring context: 194B
=> CACHED [server stage-0 2/5] WORKDIR /app
=> CACHED [server stage-0 3/5] RUN adduser     --disabled-password     --gecos ""     --home "/nonexistent"     --shell "/sbin/nolo
=> CACHED [server stage-0 4/5] RUN --mount=type=cache,target=/root/.cache/pip     --mount=type=bind,source=requirements.txt,target=
=> CACHED [server stage-0 5/5] COPY . .
=> [server] exporting to image
=> => exporting layers
=> => writing image sha256:b77c7c31474238cacb65a212bf106b413f1a9880c221b9109830618c193fbfa0
=> => naming to docker.io/library/docker-genai-sample-server
=> [server] resolving provenance for metadata file
+] Running 3/0
✔ Container docker-genai-sample-ollama-pull-1   Created
✔ Container docker-genai-sample-database-1      Running
✔ Container docker-genai-sample-ollama-1        Recreated
Attaching to database-1, ollama-1, ollama-pull-1, server-1
ollama-1       | 2024/11/26 09:37:37 routes.go:1197: INFO server config env="map[CUDA_VISIBLE_DEVICES: GPU_DEVICE_ORDINAL: HIP_VISIB
ON: HTTPS_PROXY: HTTP_PROXY: NO_PROXY: OLLAMA_DEBUG:false OLLAMA_FLASH_ATTENTION:false OLLAMA_GPU_OVERHEAD:0 OLLAMA_HOST:http://0.0
OLLAMA_KEEP_ALIVE:5m0s OLLAMA_LLM_LIBRARY: OLLAMA_LOAD_TIMEOUT:5m0s OLLAMA_MAX_LOADED_MODELS:0 OLLAMA_MAX_QUEUE:512 OLLAMA_MODELS:/
ER_CACHE:false OLLAMA_NOHISTORY:false OLLAMA_NOPRUNE:false OLLAMA_NUM_PARALLEL:0 OLLAMA_ORIGINS:[http://localhost https://localhost
host:* http://127.0.0.1 https://127.0.0.1 http://127.0.0.1:* https://127.0.0.1:* http://0.0.0.0 https://0.0.0.0 http://0.0.0.0:* htt
uri://* vscode-webview://*] OLLAMA_SCHED_SPREAD:false OLLAMA_TMPDIR: ROCR_VISIBLE_DEVICES: http_proxy: https_proxy: no_proxy:]"
llama-1        | time=2024-11-26T09:37:37.882Z level=INFO source=images.go:753 msg="total blobs: 0"
llama-1        | time=2024-11-26T09:37:37.882Z level=INFO source=images.go:760 msg="total unused blobs removed: 0"
llama-1        | time=2024-11-26T09:37:37.882Z level=INFO source=routes.go:1248 msg="Listening on [::]:11434 (version 0.4.5)"
llama-1        | time=2024-11-26T09:37:37.883Z level=INFO source=common.go:49 msg="Dynamic LLM libraries" runners="[cuda_v12 cpu cuc
"
llama-1        | time=2024-11-26T09:37:37.883Z level=INFO source=gpu.go:221 msg="looking for compatible GPUs"
llama-1        | time=2024-11-26T09:37:37.884Z level=INFO source=gpu.go:386 msg="no compatible GPUs were discovered"
llama-1        | time=2024-11-26T09:37:37.884Z level=INFO source=types.go:123 msg="inference compute" id=0 library=cpu variant="no v
r=0.0 name="" total="3.8 GiB" available="3.0 GiB"
llama-pull-1 | pulling ollama model gpt-3.5 using http://host.docker.internal:11434
llama-pull-1 | OLLAMA model only pulled if both LLM and OLLAMA_BASE_URL are set and the LLM model is not gpt
llama-pull-1 exited with code 0
erver-1        |
erver-1        | Collecting usage statistics. To deactivate, set browser.gatherUsageStats to false.
erver-1        |
erver-1        |
erver-1        |   You can now view your Streamlit app in your browser.
erver-1        |
erver-1        |   URL: http://0.0.0.0:8000
erver-1        |
erver-1        | /usr/local/lib/python3.11/site-packages/transformers/utils/generic.py:441: FutureWarning: `torch.utils._pytree._reg
```
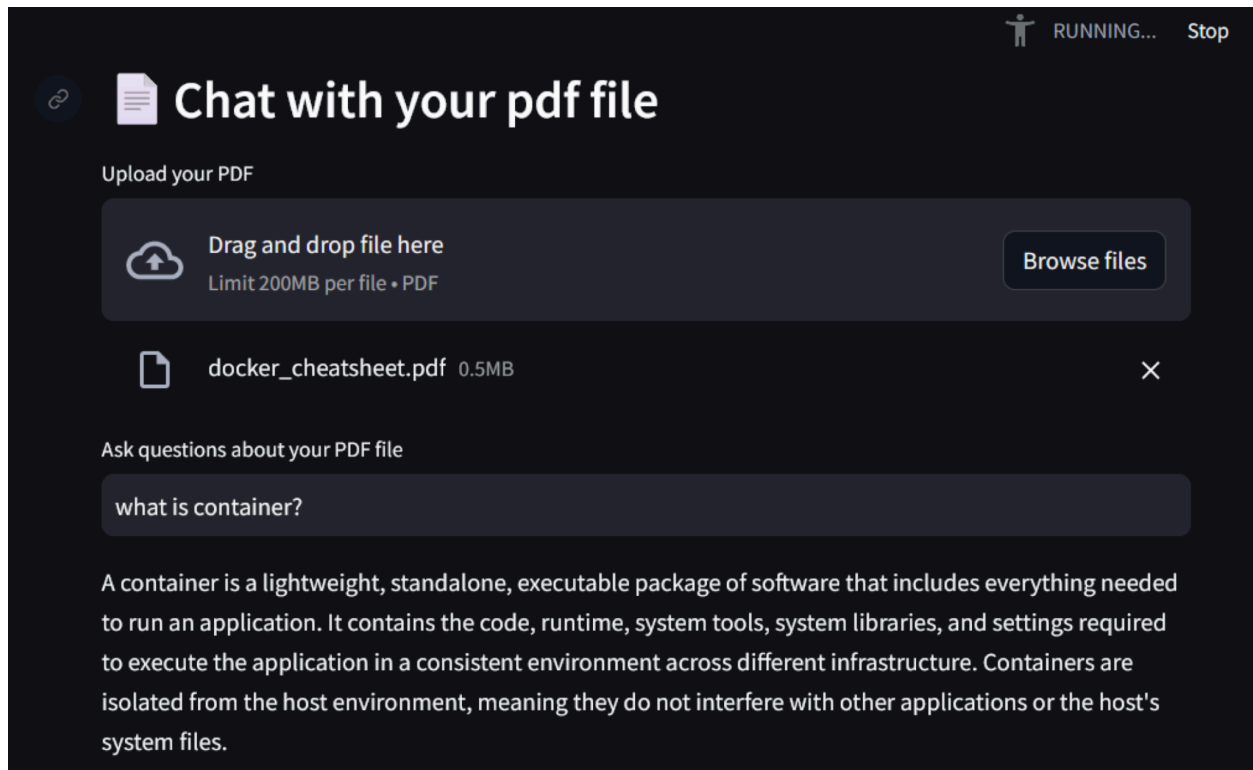
- Wait until everything is built and service is started.

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.

Week 11 Homework 2: Project: GenAI - Develop your containerized app



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

========================================================================

**Step 3: Link to GitHub**