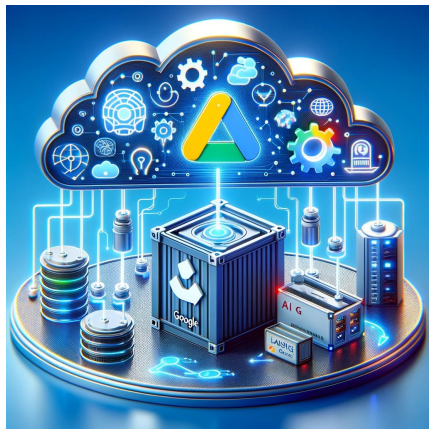# Pulling Google Drive Data into RAG with LangServe and Docker

Deploying a RAG toolchain using LangServe in your local Docker environment

# Purpose

- Learn to leverage LangChain Templates and Docker to build scalable microservices for your chains with ease

- See how you can interact with LangServe to create FastAPI endpoints for dynamically create and invoke your chains
- Streamlit as frontend to interact with your chains
- Docker to containerize and deploy your AI microservices

- See the app in action performing RAG on your own data

- Try it out and create your own with your preferred document loader!

# Use Cases for RAG on your own Data

- Retrieving engineering standards / procedures
- Querying meeting transcripts
- Getting insights on project info that would otherwise be overlooked

# Lab Deep Dive
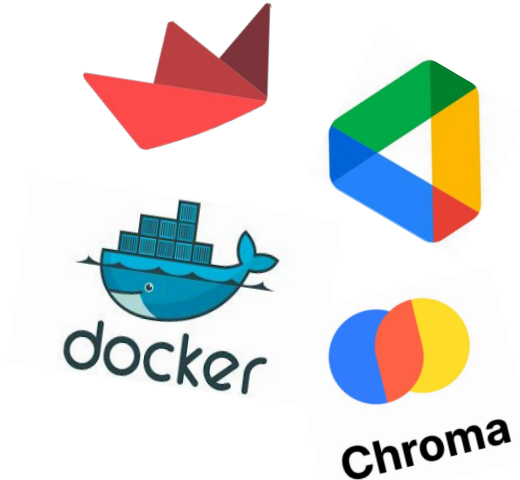
# Technologies Used

**Backend**

- LangChain rag-chroma template
- Fast API + LangServe
- Open AI - ChatOpenAI + OpenAIEmbeddings
- Chroma
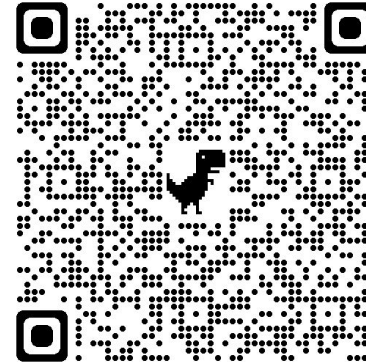- Google Drive API

**Frontend**

- Streamlit - link to tutorial

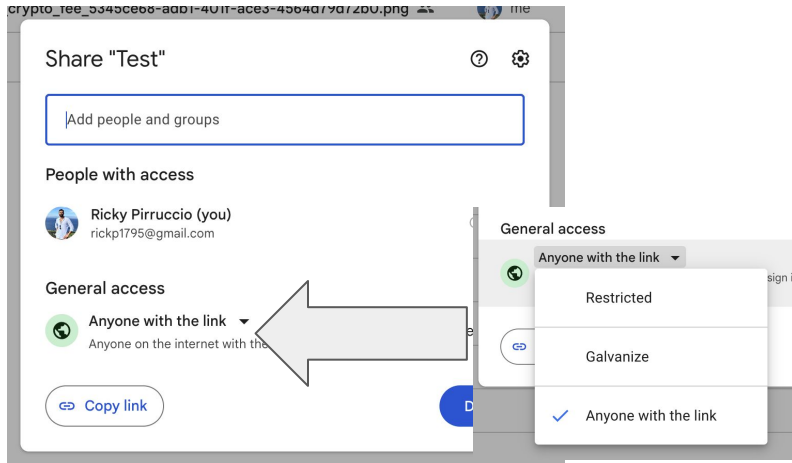**Deployment and Development Pipeline**

- Docker - link to tutorial

# Google Credentials Prerequisites

- You must create a google service account before you can access the Google Drive API on your own folders
- Karim Lalani made a great tutorial on this [here](#)
- Important: Set your folder sharing permissions to "Anyone with the link"

# Google Credentials Prerequisites

- After creating a service account, you'll receive a key file with credentials for authentication. Rename it keys.json and save it in this path ~/.credentials/keys.json

# Prerequisites

**Clone the repo**

- git clone https://github.com/colinmcnamara/austin_langchain
- cd austin_langchain/labs/LangChain_103/rag_chroma_from_google_drive

**Installs**

- [Docker Desktop](#)

**Environment Variables**

- `GOOGLE_APPLICATION_CREDENTIALS:` Set this variable to the path where you saved the <mark>keys.json</mark> file earlier. Make sure to call out the full path or you will get an error
- `OPENAI_API_KEY`

# Live Demo

# Docker setup

Chain Microservice
(Backend)

Streamlit Microservice
(Frontend)

Compose

```dockerfile
Dockerfile > ...
1   FROM python:3.11-slim
2
3   RUN pip install poetry==1.6.1
4
5   RUN pip install google-api-python-client google-auth-oauthlib && \
6       pip install pypdf2
7
8   RUN poetry config virtualenvs.create false
9
10  ENV GOOGLE_APPLICATION_CREDENTIALS=/app/.credentials/keys.json
11
12  WORKDIR /code
13
14  COPY ./pyproject.toml ./langchain_template_README.md ./poetry.lock* ./
15
16  COPY ./packages ./packages
17
18  RUN poetry install --no-interaction --no-ansi --no-root
19
20  COPY ./app ./app
21
22  RUN poetry install --no-interaction --no-ansi
23
24  EXPOSE 8000
25
26  CMD exec uvicorn app.server:app --host 0.0.0.0 --port 8000
27
```

```dockerfile
Dockerfile.streamlit > ...
1   # Use an official Python runtime as a parent image
2   FROM python:3.11-slim
3
4   # Set the working directory in the container
5   WORKDIR /usr/src/app
6
7   # Install Streamlit
8   RUN pip install streamlit
9
10  # Copy the Streamlit app file into the container
11  COPY streamlit_chat.py .
12
13  # Make port 8501 available to the world outside this container
14  EXPOSE 8501
15
16  # Run Streamlit when the container launches
17  CMD ["streamlit", "run", "streamlit_chat.py"]
```

🐳 docker-compose.yml
🐳 Dockerfile
🐳 Dockerfile.streamlit

```yaml
docker-compose.yml
1   version: '3.8'
2
3   services:
4     api:
5       build:
6         context: .
7         dockerfile: Dockerfile
8       environment:
9         - OPENAI_API_KEY=${OPENAI_API_KEY}
10      volumes:
11        - ${GOOGLE_APPLICATION_CREDENTIALS}:/app/.credentials/keys.json
12      ports:
13        - "8000:8000"
14
15    frontend:
16      build:
17        context: .
18        dockerfile: Dockerfile.streamlit
19      depends_on:
20        - api
21      ports:
22        - "8501:8501"
```

# Basic Backend Setup

Retrieves a chain

Adds a new chain

```python
@app.get("/")
async def redirect_root_to_docs():
    return RedirectResponse("/docs")


@app.get("/list-folder-routes")
def list_folder_routes():
    routes = set()
    for route in app.routes:
        if isinstance(route, APIRoute):
            if route.name == "invoke" and route.path.startswith(
                "/folders/{folder_id}/"
            ):
                routes.add(route.path)
    return list(routes)


@app.post("/initialize-chain")
async def initialize_chain_endpoint(
    folder_id: str = Body(..., embed=True), name: str = Body(..., embed=True)
):
    try:
        chain = create_chain(folder_id)
        new_path = f"/folders/{folder_id}/{name}"
        add_routes(app, chain, path=new_path)
        return {
            "message": f"Chain initialized successfully at {new_path}",
            "path": f"{new_path}/invoke",
        }
    except Exception as e:
        raise HTTPException(status_code=500, detail=str(e))


if __name__ == "__main__":
    import uvicorn

    uvicorn.run(app, host="0.0.0.0", port=8000)
```

rap chain creation in a function. Use to create FastAPI chain endpoint

Doc loader

Text splitter

Vector store layer

Prompt

Chain

Model

```python
def create_chain(folder_id):


    loader = GoogleDriveLoader(
        folder_id=folder_id,
        recursive=False,
        # we need to use service_account_key to get google credentials because we're using do
        issues/8755
        file_types=["document", "sheet", "pdf"],
        service_account_key=os.environ["GOOGLE_APPLICATION_CREDENTIALS"],
    )
    data = loader.load()


    # Split
    text_splitter = RecursiveCharacterTextSplitter(
        chunk_size=1000, chunk_overlap=200, separators=[" ", ",", "\n"]
    )
    all_splits = text_splitter.split_documents(data)


    # Add to vectorDB
    vectorstore = Chroma.from_documents(
        documents=all_splits,
        collection_name="rag-chroma",
        embedding=OpenAIEmbeddings(),
    )
    retriever = vectorstore.as_retriever()

    # RAG prompt
    template = """
    Answer the question based only on the following context:
    {context}

    Question: {question}
    """
    prompt = ChatPromptTemplate.from_template(template)

    # LLM
    model = ChatOpenAI()

    # RAG chain
    chain = (
        RunnableParallel({"context": retriever, "question": RunnablePassthrough()})
        | prompt
        | model
        | StrOutputParser()
    )

    # Add typing for input
    class Question(BaseModel):
        __root__: str


    chain = chain.with_types(input_type=Question)

    return chain
```

# Basic Frontend Setup

Handles selecting and displaying chains in the sidebar

Define vars

```python
import requests
import streamlit as st

API_BASE = "http://api:8000"
state = st.session_state
```

Initializes state

Handles saving / retrieving chat history

```python
def handle_folder_configs():
    if state.folder_routes == []:
        get_folder_routes()

    with st.sidebar:
        st.header("Folder Configurations")
        folder_id = st.text_input("Enter Google Drive Folder ID")
        config_name = st.text_input("Enter a name for this configuration")
        if st.button(
            "Add New Configuration",
        ):
            add_new_folder_config(config_name, folder_id)

        if state.folder_routes:
            st.sidebar.header("Select Drive to Chat With")

            for path in state.folder_routes:
                if st.button(path):
                    select_config(path)
```

```python
def initialize_state():
    if "chat_histories" not in state:
        state.chat_histories = {}
    if "selected_path" not in state:
        state.selected_path = ""
    if "folder_routes" not in state:
        state.folder_routes = []
```

```python
class ChatManager:
    def add_message(self, role, content):
        if state.selected_path not in state.chat_histories:
            state.chat_histories[state.selected_path] = []

        message = {
            "role": role,
            "content": content,
        }
        state.chat_histories[state.selected_path].append(message)

        with st.chat_message(role):
            st.write(content)

    def display_chat_history(self):
        # Display chat messages from history for the selected path
        if state.selected_path in state.chat_histories:
            for message in state.chat_histories[state.selected_path]:
                with st.chat_message(message["role"]):
                    st.write(message["content"])
```

Our condensed app

```python
st.title("Chat With Google Drive Files – Rag Chroma")

initialize_state()
chat_manager = ChatManager()

handle_folder_configs()

if state.selected_path:
    chat_manager.display_chat_history()

    if prompt := st.chat_input("What is your query?"):
        chat_manager.add_message("user", prompt)

        response = get_response_from_llm(prompt, state.selected_path)

        chat_manager.add_message("assistant", response)
```

# Connecting Backend with Frontend

Query a chain

Get list of chain routes
(GET request endpoint)

Create new chain / chain
endpoint
(POST request endpoint)

```python
@st.cache_data
def get_response_from_llm(query, path):
    payload = {"input": query, "config": {}, "kwargs": {}}
    res = requests.post(f"{API_BASE}{path}", json=payload)
    if res.status_code == 200:
        return res.json()["output"]
    else:
        return f"Error: Received status code {res.status_code}"


def get_folder_routes():
    response = requests.get(f"{API_BASE}/list-folder-routes")
    if response.status_code == 200:
        folder_routes = response.json()
        state.folder_routes = folder_routes
    else:
        st.error(f"{response.status_code} Error: Failed to retrieve folder routes")


def add_new_folder_config(name, folder_id):
    response = requests.post(
        f"{API_BASE}/initialize-chain", json={"folder_id": folder_id, "name": name}
    )
    if response.status_code == 200:
        path = response.json().get("path")
        state.folder_routes.append(path)
        select_config(path)
    else:
        error_detail = response.json().get("detail")

        if "File not found" in error_detail:
            user_friendly_error = "Invalid folder ID, please try again."
        else:
            user_friendly_error = error_detail

        st.error(f"{response.status_code} Error: {user_friendly_error}")
```