

Chat with your PDF - Gen AI App - With Amazon Bedrock, RAG, S3, Langchain and Streamlit

we will build a CHATBOT like application with AWS Amazon Bedrock, docker, python, Langchain, and Streamlit. We will use Retrieval-Augmented generation concept to provide context to the Large Language model along with user query to generate response from our Knowledgebase

I will demonstrate the following: -

Architecture of the applications - Build 2 applications (ADMIN and USER) and create DOCKER images –

ADMIN Application: - Build Admin Web application where AdminUser can upload the pdf. - The PDF text is split into chunks - Using the Amazon Titan Embedding Model, create the vector representation of the chunks - Using FAISS, save the vector index locally - Upload the index to Amazon S3 bucket (You can use other vector stores like OpenSearch, Pinecone, PgVector etc., but for this demo, I chose cost effective S3) –

USER Application: - Build User Web application where users can query / chat with the pdf. - At the application start, download the index files from S3 to build local FAISS index (vector store) - Langchain's RetrievalQA, does the following: - Convert the User's query to vector embedding using Amazon Titan Embedding Model (Make sure to use the same model that was used for creating the chunk's embedding on the Admin side) - Do similarity search to the FAISS index and retrieve 5 relevant documents pertaining to the user query to build the context - Using Prompt template, provide the question and context to the Large Language Model. We are using Claude model from Anthropic. - Display the LLM's response to the user.

PDF .. Loaded in this folder- <https://www.versusarthritis.org/media/22726/what-is-arthritis-information-booklet.pdf>

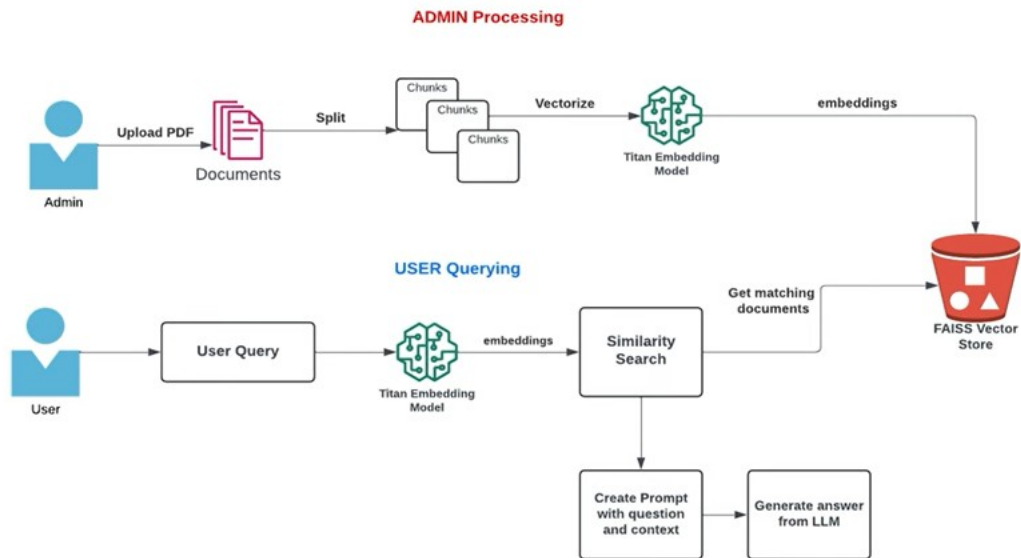
Amazon S3

youtube-rag-test-gj

Objects (2)

Name	Type	Last modified	Size	Storage class
my_faiss.faiss	faiss	May 2, 2024, 13:10:44 (UTC-07:00)	252.0 KB	Standard
my_faiss.pkl	pkl	May 2, 2024, 13:10:45 (UTC-07:00)	38.1 KB	Standard

What will we build?



What is RAG

- ✓ RAG stands for **Retrieval-Augmented Generation**
- ✓ Process of optimizing the output of Large Language Model
- ✓ LLM uses knowledge base outside its training data source
- ✓ To query LLM: Provide question and context in the Prompt
- ✓ Depending on use-case it can be Cost Effective compared to Fine Tuning

Build Admin Site

- ✓ Create S3 Bucket
- ✓ Create requirements.txt
- ✓ Create Python App
- ✓ Create Dockerfile and build docker image
- ✓ Access the application from Browser
- ✓ Upload a PDF
- ✓ Confirm the FAISS vector index files are uploaded to S3

Build User Site

- ✓ Create requirements.txt
- ✓ Create Python App
- ✓ Create Dockerfile and build docker image
- ✓ Access the application from Browser
- ✓ Ask a question
- ✓ Check the response

Go to aws console - s3- create S3 bucket

Create bucket [Info](#)

Buckets are containers for data stored in S3.

General configuration

AWS Region

US West (Oregon) us-west-2

Bucket type [Info](#)

☒ **General purpose**

Recommended for most use cases and access patterns. General purpose buckets are the original S3 bucket type. They allow a mix of storage classes that redundantly store objects across multiple Availability Zones.

☐ **Directory - New**

Recommended for low-latency use cases. These buckets use only the S3 Express One Zone storage class, which provides faster processing of data within a single Availability Zone.

Bucket name [Info](#)

myawsbucket

Bucket name must be unique within the global namespace and follow the bucket naming rules. [See rules for bucket naming](#)

Copy settings from existing bucket - optional

Only the bucket settings in the following configuration are copied.

[Choose bucket](#)

Format: s3://bucket/prefix

Object Ownership [Info](#)

Control ownership of objects written to this bucket from other AWS accounts and the use of access control lists (ACLs). Object ownership determines who can specify access to objects.

☒ **ACLs disabled (recommended)**

All objects in this bucket are owned by this account.

Access to this bucket and its objects is specified using only policies.

☐ **ACLs enabled**

Objects in this bucket can be owned by other AWS accounts.

Access to this bucket and its objects can be specified using ACLs.

Bucket name [Info](#)

youtube-demo-bedrock-chat-pdf

Bucket name must be unique within the global namespace and follow the bucket naming rules. [See rules for bucket naming](#)

Copy settings from existing bucket - optional

Only the bucket settings in the following configuration are copied.

[Choose bucket](#)

Format: s3://bucket/prefix

Object Ownership [Info](#)

Control ownership of objects written to this bucket from other AWS accounts and the use of access control lists (ACLs). Object ownership determines who can specify access to objects.

☒ **ACLs disabled (recommended)**

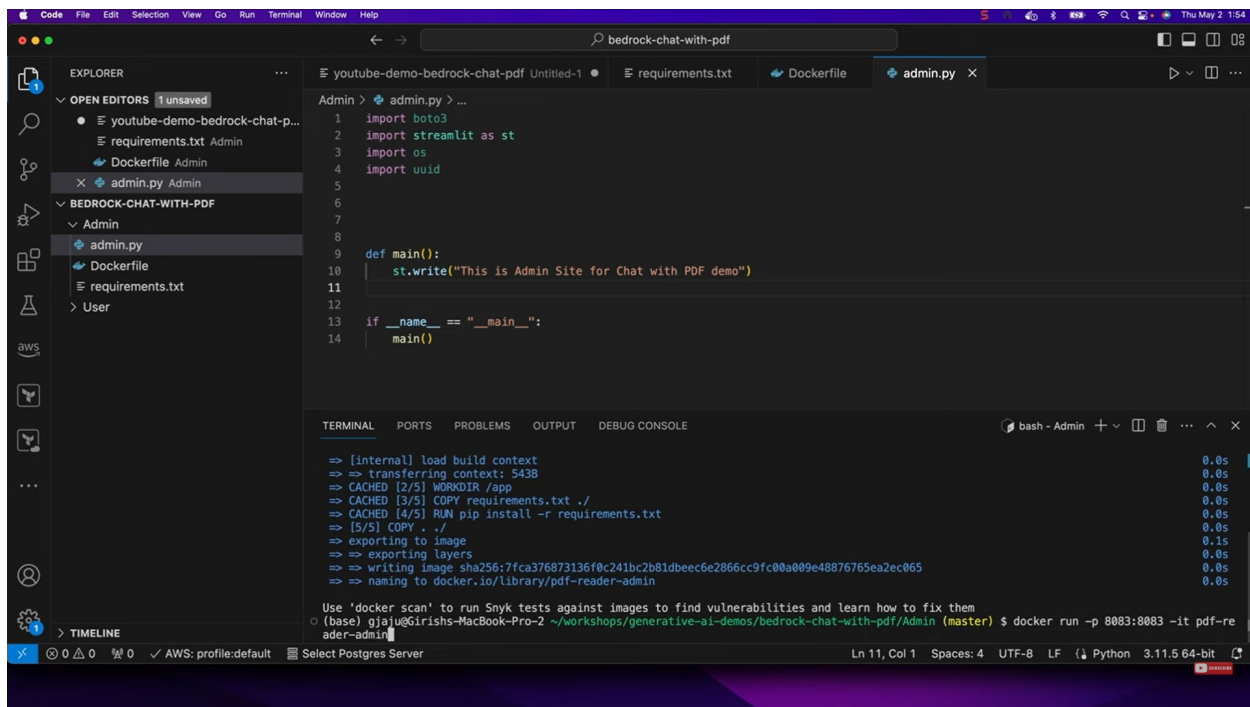
All objects in this bucket are owned by this account. Access to this bucket and its objects is specified using only policies.

☐ **ACLs enabled**

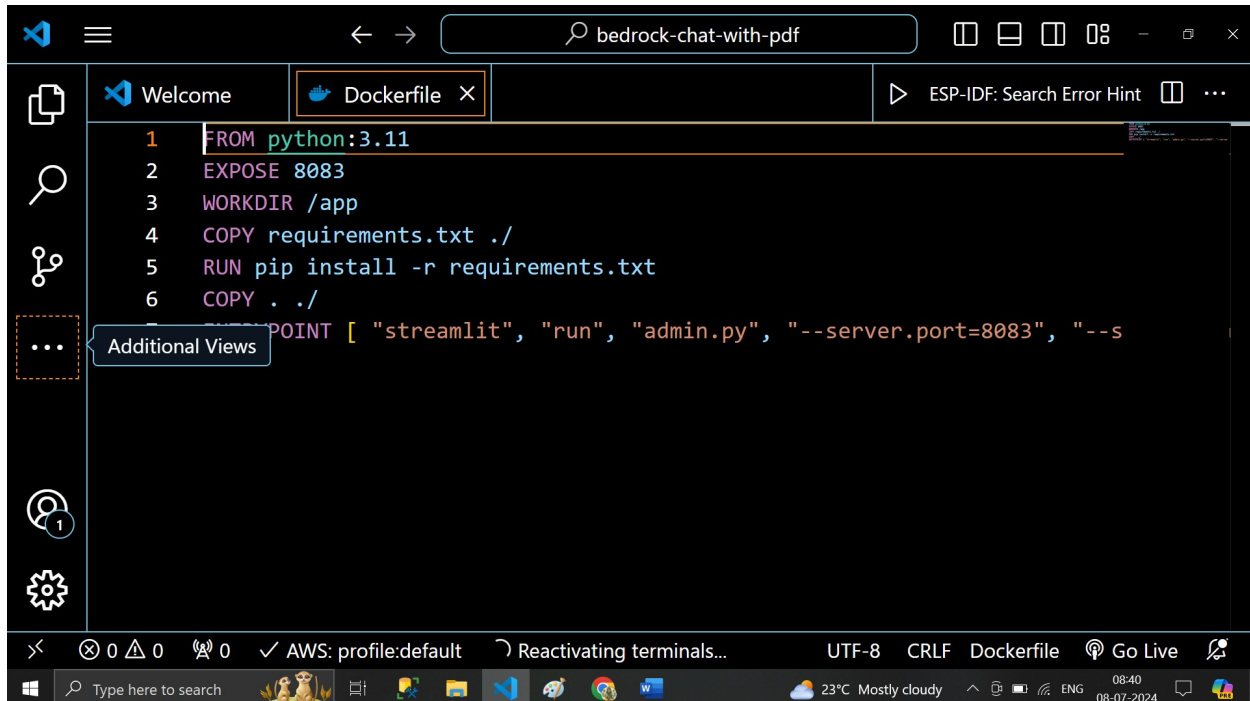
Objects in this bucket can be owned by other AWS accounts. Access to this bucket and its objects can be specified using ACLs.

Object Ownership

Bucket owner enforced



docker file



```
File Edit Selection View Go ... bedrock-chat-with-pdf
admin.py
1 import boto3
2 import streamlit as st
3 import os
4 import uuid
5
6 ## s3 client
7 s3_client = boto3.client("s3")
8 BUCKET_NAME = os.getenv("BUCKET_NAME")
9
10 ## Bedrock
11 from langchain_community.embeddings import BedrockEmbeddings
12
13 ## Text Splitter
14 from langchain.text_splitter import RecursiveCharacterTextSplitter
15
16 ## Pdf Loader
17 from langchain_community.document_loaders import PyPDFLoader
18
19 ## import FAISS
20 from langchain_community.vectorstores import FAISS
21
22 bedrock_client = boto3.client(service_name="bedrock-runtime")
23 bedrock_embeddings = BedrockEmbeddings(model_id="amazon.titan-embed-text-v1", client=bedrock_client)
24
25 def get_unique_id():
26     return str(uuid.uuid4())
27
28
29 ## Split the pages / text into chunks
30 def split_text(pages, chunk_size, chunk_overlap):
31     text_splitter = RecursiveCharacterTextSplitter(chunk_size=chunk_size, chunk_overlap=chunk_overlap)
32     docs = text_splitter.split_documents(pages)
33     return docs
34
35 ## create vector store
36 def create_vector_store(request_id, documents):
37     vectorstore_faiss = FAISS.from_documents(documents, bedrock_embeddings)
38     file_name = f"{request_id}.bin"
```

Admin -side .. AWS Bedrock

Amazon Bedrock

The new Titan Embeddings G1 – Text v1.2 can intake up to 8k tokens and outputs a vector of 1536 dimensions. The model also works in 25+ different languages. The model is optimized for text retrieval tasks but can also perform additional tasks such as semantic similarity and clustering. Titan Embeddings G1 – Text v1.2 also supports long documents, however, for retrieval tasks it is recommended to segment documents into logical segments (such as paragraphs or sections).

Supported formats
text retrieval, semantic similarity, clustering
[View example notebook](#)

Model attributes
Output vector size = 1,536

Model version
v1.2

Max tokens
8k

Languages
English, Arabic, Chinese (Sim.), French, German, Hindi, Japanese, Spanish, Czech, Filipino, Hebrew, Italian, Korean, Portuguese, Russian, Swedish, Turkish, Chinese (trad), Dutch, Kannada, Malayalam, Marathi, Polish, Tamil, Telugu and others.

End user license agreement (EULA)
[View EULA](#)

Pricing
[View pricing](#)

Resource
[Bedrock documentation on embeddings](#)

API request

```
{
  "modelId": "amazon.titan-embed-text-v1",
  "contentType": "application/json",
  "accept": "*/*",
  "body": "{\\\"inputText\\\":\\\"this is where you place your input text\\\"}"
}
```

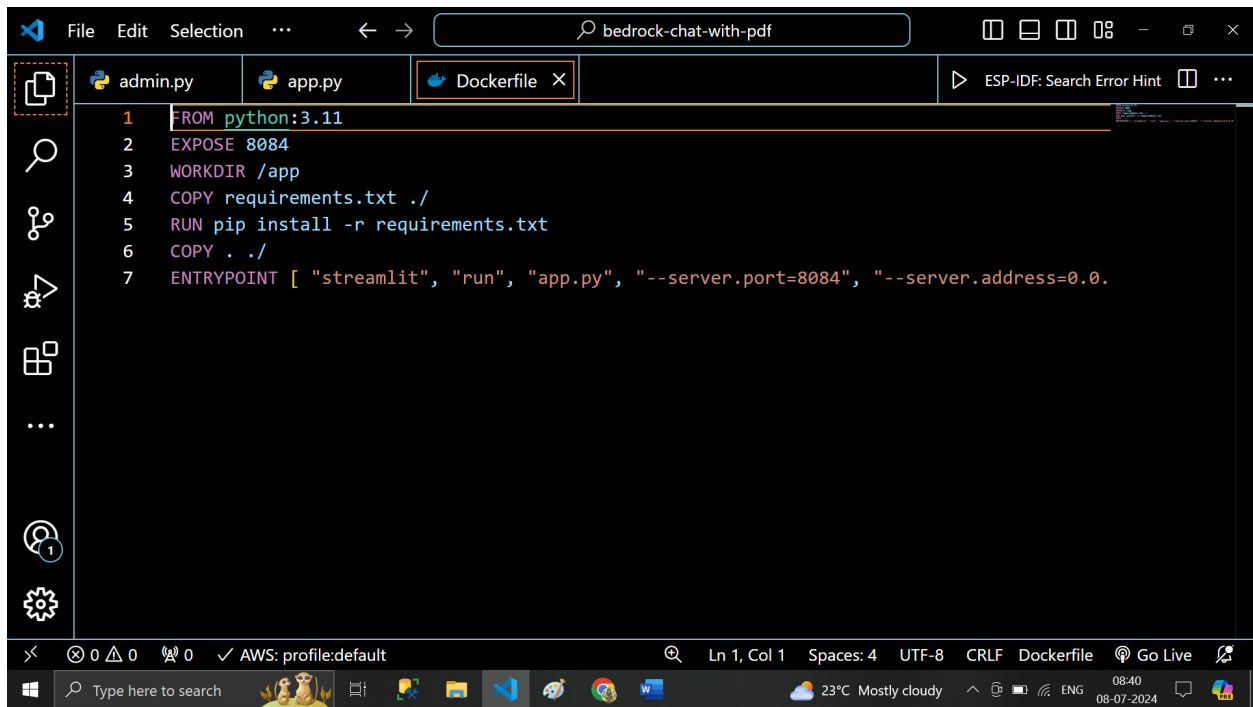
Resources

Model access [new](#)

Settings

User guide [📖](#)

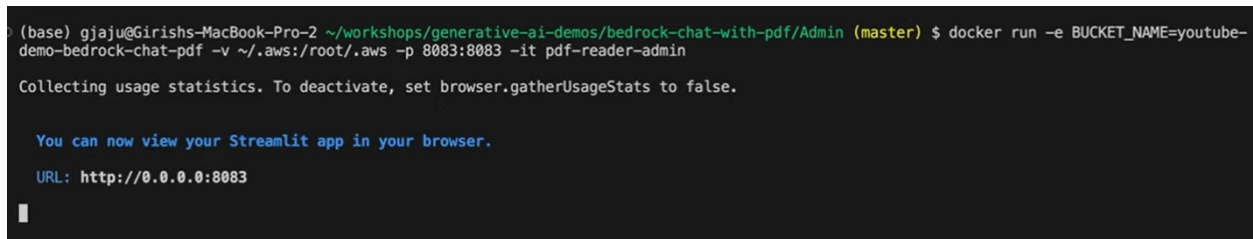
© 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences



The image shows a Visual Studio Code editor window titled "bedrock-chat-with-pdf". The "Dockerfile" tab is active, displaying the following content:

```
1 FROM python:3.11
2 EXPOSE 8084
3 WORKDIR /app
4 COPY requirements.txt ./
5 RUN pip install -r requirements.txt
6 COPY . ./
7 ENTRYPOINT [ "streamlit", "run", "app.py", "--server.port=8084", "--server.address=0.0.0.0" ]
```

The status bar at the bottom indicates the file is at "Ln 1, Col 1" with "Spaces: 4", "UTF-8" encoding, and "CRLF" line endings. The "Go Live" button is also visible.



The terminal shows the execution of the following command:

```
(base) gjaju@Girishs-MacBook-Pro-2 ~/workshops/generative-ai-demos/bedrock-chat-with-pdf/Admin (master) $ docker run -e BUCKET_NAME=youtube-demo-bedrock-chat-pdf -v ~/.aws:/root/.aws -p 8083:8083 -it pdf-reader-admin
```

The output indicates that the application is running successfully and provides the URL to view it:

```
Collecting usage statistics. To deactivate, set browser.gatherUsageStats to false.

You can now view your Streamlit app in your browser.

URL: http://0.0.0.0:8083
```

User- Data- app.py file


```
File Edit Selection View Go ... bedrock-chat-with-pdf
app.py x
1 import boto3
2 import streamlit as st
3 import os
4 import uuid
5
6 ## s3_client
7 s3_client = boto3.client("s3")
8 BUCKET_NAME = os.getenv("BUCKET_NAME")
9
10 ## Bedrock
11 from langchain_community.embeddings import BedrockEmbeddings
12 from langchain.llms.bedrock import Bedrock
13
14 ## prompt and chain
15 from langchain.prompts import PromptTemplate
16 from langchain.chains import RetrievalQA
17
18 ## Text Splitter
19 from langchain.text_splitter import RecursiveCharacterTextSplitter
20
21 ## Pdf Loader
22 from langchain_community.document_loaders import PyPDFLoader
23
24 ## import FAISS
25 from langchain_community.vectorstores import FAISS
```

Code Editor Interface showing a project named "youtube-demo-bedrock-chat-pdf".

EXPLORER:

- OPEN EDITORS (1 unsaved)
 - youtube-demo-bedrock-chat-pdf
 - requirements.txt Admin
 - Dockerfile Admin
 - admin.py Admin
 - requirements.txt User
 - Dockerfile User
 - app.py User
- BEDROCK-CHAT-WITH-PDF
 - Admin
 - admin.py
 - Dockerfile
 - requirements.txt
 - User
 - app.py
 - Dockerfile
 - requirements.txt

EDITOR: youtube-demo-bedrock-chat-pdf

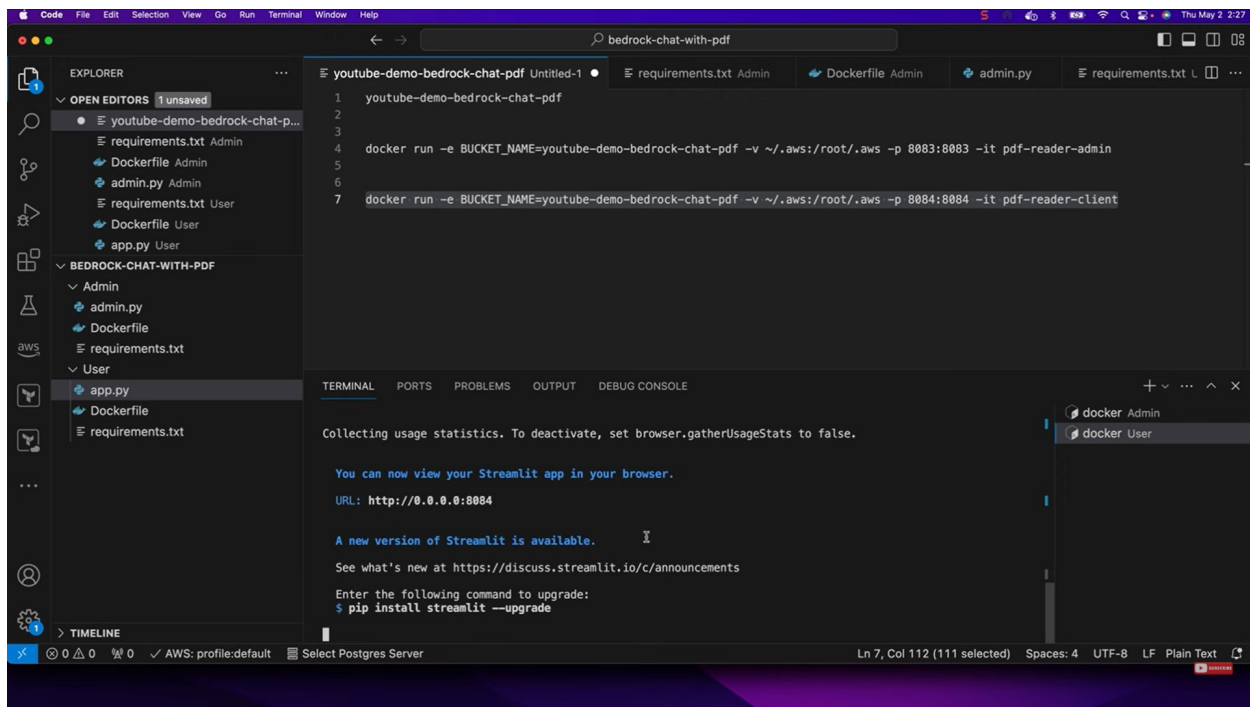
```
1 youtube-demo-bedrock-chat-pdf
2
3
4 docker run -e BUCKET_NAME=youtube-demo-bedrock-chat-pdf -v ~/.aws:/root/.aws -p 8083:8083 -it pdf-reader-admin
5
6
7 docker run -e BUCKET_NAME=youtube-demo-bedrock-chat-pdf -v ~/.aws:/root/.aws -p 8084:8084 -it pdf-reader-client
```

TERMINAL:

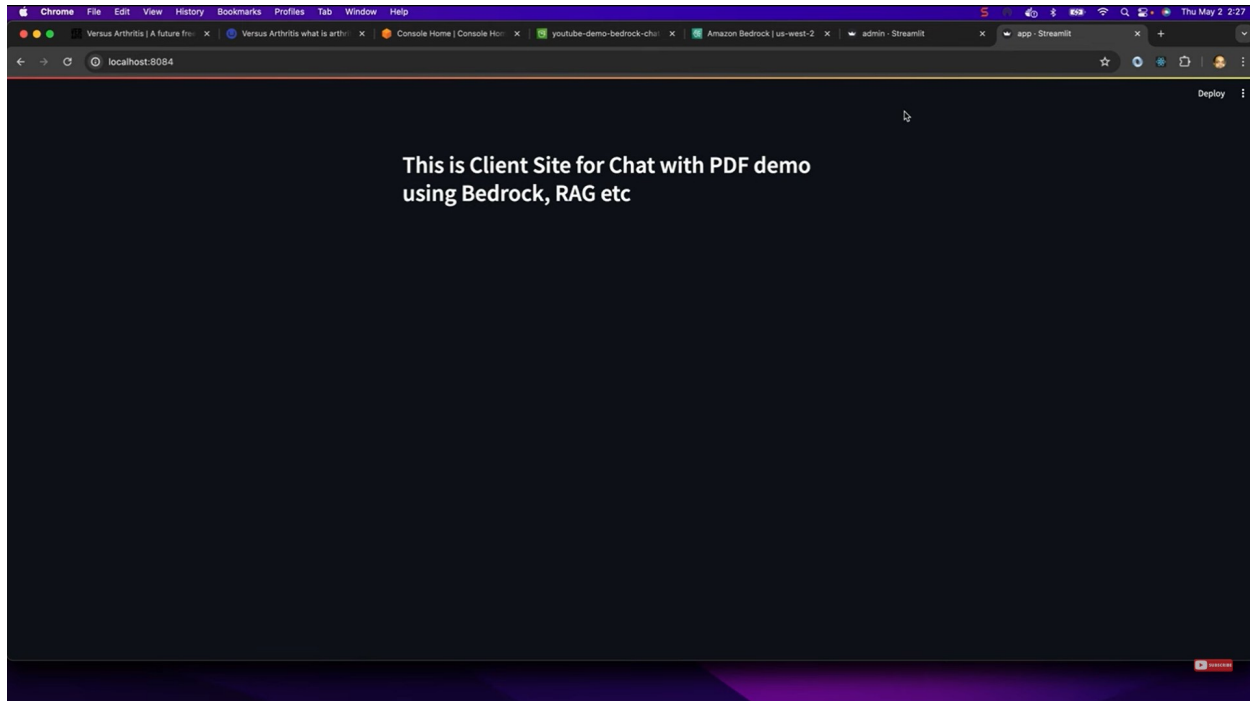
```
=> [internal] load build context                                0.0s
=> => transferring context: 1.18kB                               0.0s
=> [1/5] FROM docker.io/library/python:3.11@sha256:e453eb723bc8ecac7a797498f9a5915d13e567628d48cd356 0.0s
=> CACHED [2/5] WORKDIR /app                                    0.0s
=> CACHED [3/5] COPY requirements.txt ./                          0.0s
=> CACHED [4/5] RUN pip install -r requirements.txt              0.0s
=> [5/5] COPY . ./                                              0.0s
=> exporting to image                                           0.1s
=> exporting layers                                             0.0s
=> writing image sha256:86edd29077c4697bad8d5dc0757eafe80fe16ab61b0322332c02b36d9cbee7 0.0s
=> naming to docker.io/library/pdf-read-client                 0.0s

Use 'docker scan' to run Snyk tests against images to find vulnerabilities and learn how to fix them
(base) gjaugGirishs-MacBook-Pro-2 ~/workshops/generative-ai-demos/bedrock-chat-with-pdf/User (master) $ dock
er build -t pdf-reader-client .
[*] Building 0.0s (0/1)
```

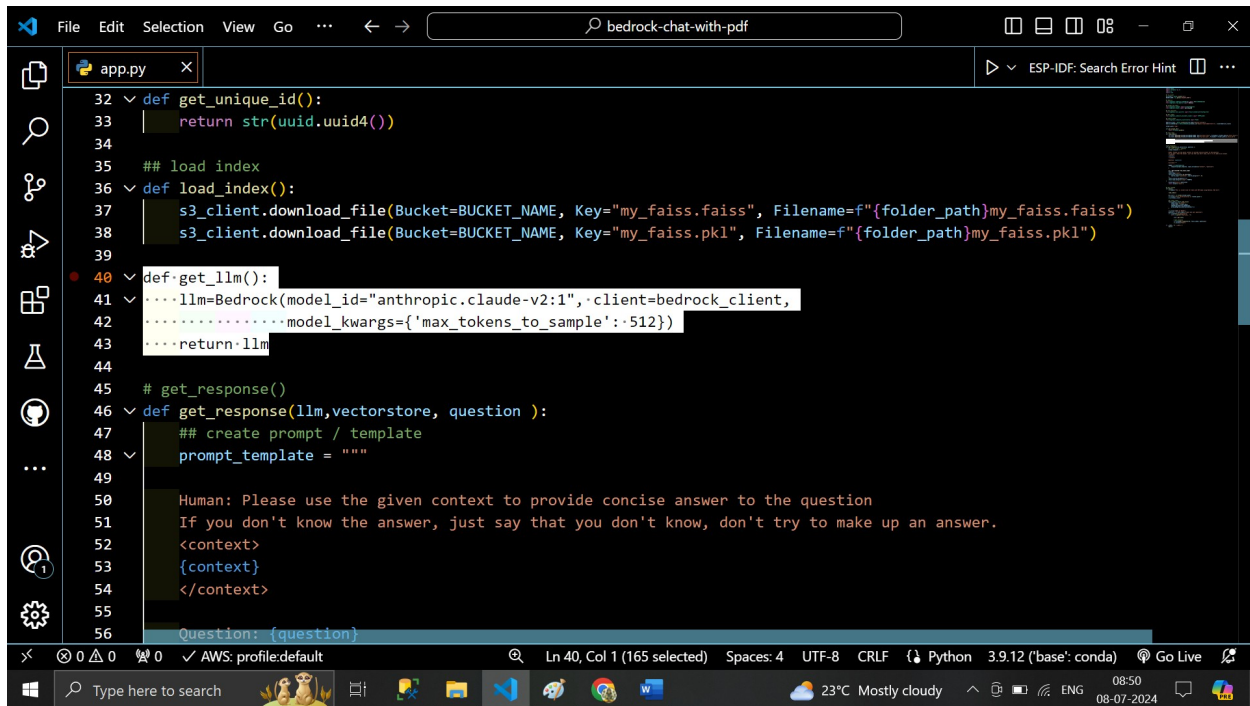
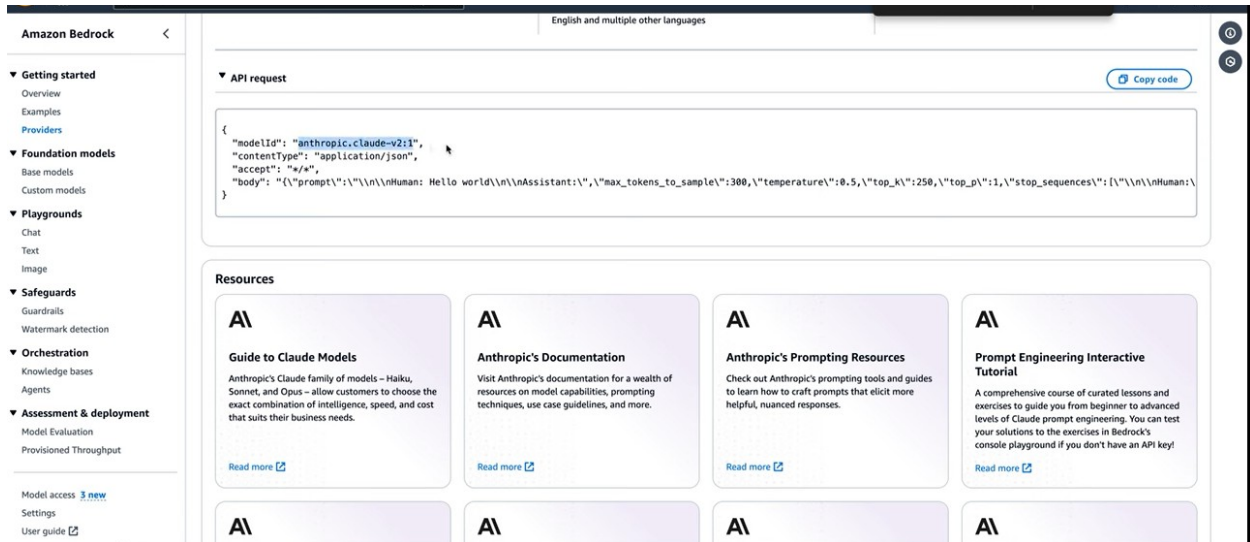
STATUS BAR: AWS profile: default, Ln 7, Col 112 (111 selected), Spaces: 4, UTF-8, LF, Plain Text



Loading app



Go to bedrock base Models



Output .. index it ready

This is Client Site for Chat with PDF demo using Bedrock, RAG etc

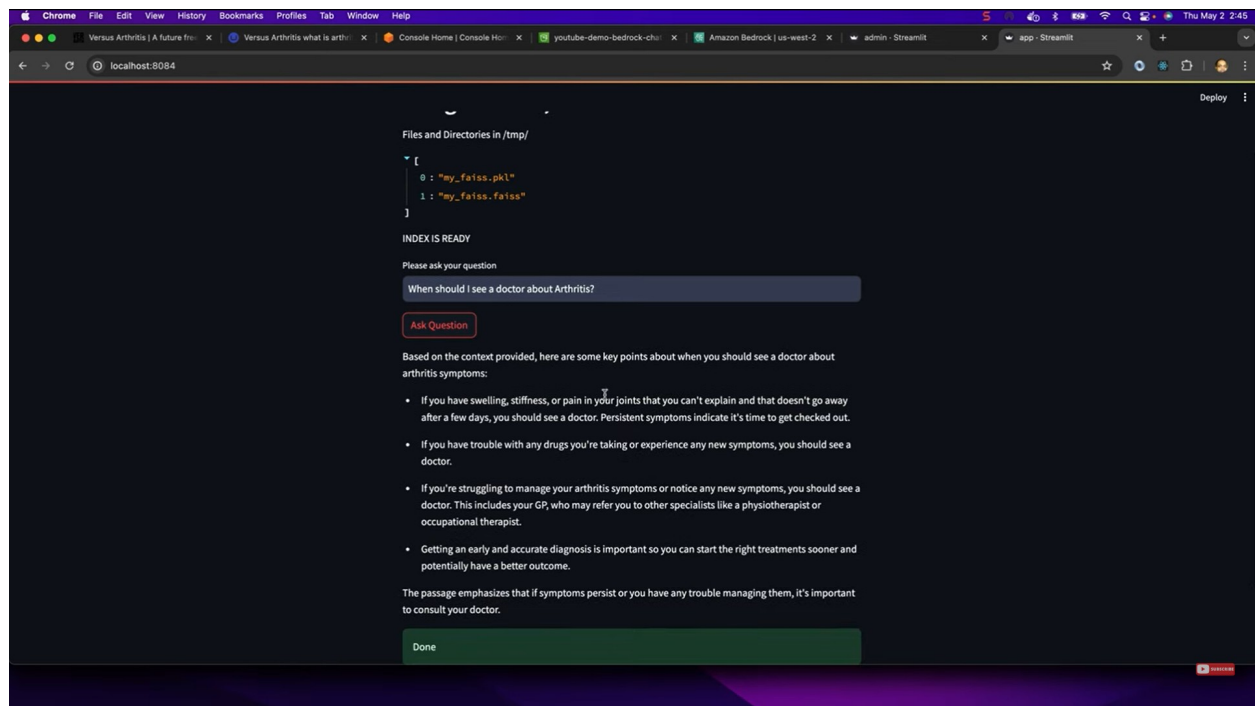
Files and Directories in /tmp/

```
[
  0 : "my_faiss.pkl"
  1 : "my_faiss.faiss"
]
```

INDEX IS READY

Please ask your question

Ask Question



Ask question Outside PDF File

INDEX IS READY

Please ask your question

Who won the cricket world cup in 2023?

Ask Question

Unfortunately, I do not have enough context in the given information to determine who won the cricket world cup in 2023. The passage discusses arthritis and provides information about Versus Arthritis, an arthritis charity. It does not mention anything about cricket or the 2023 world cup. Since there is no relevant information provided to answer the question, I don't know who won the cricket world cup in 2023.

Done