

## Source code

```
from fastapi import FastAPI, UploadFile, File, Body, status
from pydantic import BaseModel
from PyPDF2 import PdfReader
from typing import List
import faiss
import numpy as np
import sqlite3
from datetime import datetime
import os

# CONFIG
EMBEDDING_PROVIDER = "sentence_transformers"
# options: "sentence_transformers", "huggingface", "openai"

# APP
app = FastAPI(
    title="Document Q&A RAG Service",
    version="1.0.0"
)

# VECTOR DB
EMBEDDING_DIM = 384
vector_index = faiss.IndexFlatL2(EMBEDDING_DIM)
stored_chunks: List[dict] = []

# METADATA DB
conn = sqlite3.connect("metadata.db", check_same_thread=False)
cursor = conn.cursor()
cursor.execute("""
CREATE TABLE IF NOT EXISTS documents (
    id INTEGER PRIMARY KEY AUTOINCREMENT,
    filename TEXT,
    upload_time TEXT
)
""")
conn.commit()

# EMBEDDING MODELS
if EMBEDDING_PROVIDER == "sentence_transformers":
```

```

from sentence_transformers import SentenceTransformer
embedding_model = SentenceTransformer("all-MiniLM-L6-v2")

if EMBEDDING_PROVIDER == "huggingface":
    from transformers import AutoTokenizer, AutoModel
    import torch
    tokenizer = AutoTokenizer.from_pretrained("sentence-transformers/all-MiniLM-L6-v2")
    hf_model = AutoModel.from_pretrained("sentence-transformers/all-MiniLM-L6-v2")

if EMBEDDING_PROVIDER == "openai":
    from openai import OpenAI
    client = OpenAI(api_key=os.getenv("OPENAI_API_KEY"))

def get_embeddings(texts: List[str]) -> List[List[float]]:
    if EMBEDDING_PROVIDER == "sentence_transformers":

        return embedding_model.encode(texts).tolist()
    if EMBEDDING_PROVIDER == "huggingface":
        inputs = tokenizer(texts, padding=True, truncation=True, return_tensors="pt")
        with torch.no_grad():
            outputs = hf_model(**inputs)
        return outputs.last_hidden_state.mean(dim=1).tolist()
    if EMBEDDING_PROVIDER == "openai":
        response = client.embeddings.create(
            model="text-embedding-3-small",
            input=texts
        )
        return [item.embedding for item in response.data]
    raise ValueError("Invalid embedding provider")

# HELPERS
def extract_text(file: UploadFile) -> str:
    if file.filename.endswith(".pdf"):

        reader = PdfReader(file.file)
        return " ".join(page.extract_text() for page in reader.pages)
    return file.file.read().decode("utf-8")

def chunk_text(text: str, size: int = 500) -> List[str]:
    words = text.split()

```

```
return [" ".join(words[i:i + size]) for i in range(0, len(words), size)]\n\n# SCHEMAS\nclass QueryRequest(BaseModel):\n    question: str\n\nclass Source(BaseModel):\n    filename: str\n    content: str\n\nclass QueryResponse(BaseModel):\n    question: str\n    answer: str\n    sources: List[Source]\n\n# ENDPOINTS\n@app.post("/upload", status_code=status.HTTP_201_CREATED)\nasync def upload_document(file: UploadFile = File(...)):\n    text = extract_text(file)\n    chunks = chunk_text(text)\n    embeddings = get_embeddings(chunks)\n    vector_index.add(np.array(embeddings))\n    for chunk in chunks:\n        stored_chunks.append(\n            {"text": chunk,\n             "filename": file.filename\n            })\n    cursor.execute(\n        "INSERT INTO documents (filename, upload_time) VALUES (?, ?)",\n        (file.filename, datetime.now().isoformat()))\n    )\n    conn.commit()\n    return {\n        "message": "Document uploaded successfully",\n        "filename": file.filename,\n        "chunks_stored": len(chunks)\n    }\n\n@app.post("/query", response_model=QueryResponse)\nasync def query_document(payload: QueryRequest = Body(...)):\n    question_embedding = get_embeddings([payload.question])[0]
```

```
_ indices = vector_index.search(np.array([question_embedding]), k=3)
sources = []
context = ""
for i in indices[0]:
    chunk = stored_chunks[i]
    context += chunk["text"] + " "
    sources.append(
        Source(
            filename=chunk["filename"],
            content=chunk["text"]
        )
    )
answer = f"Answer based on document context:\n{context}"
return QueryResponse(
    question=payload.question,
    answer=answer,
    sources=sources
)
@app.get("/health")
def health():
    return {
        "status": "OK",
        "embedding_provider": EMBEDDING_PROVIDER
    }
```