

```
1 # docuchat_backend/app/routers/health.py
2 from __future__ import annotations
3
4 import os, inspect
5 from fastapi import APIRouter
6 from sqlalchemy import text
7
8 from ..database import engine
9 from .. import rag # only for Chroma
   ensure; NOT used for OpenAI
10 from openai import OpenAI
11 import openai as openai_pkg
12 import httpx
13
14 router = APIRouter(prefix="/health", tags
   =["health"])
15
16 def _make_openai_client() -> OpenAI:
17     api_key = os.getenv("OPENAI_API_KEY")
18     if not api_key:
19         raise RuntimeError(
20             "missing_api_key")
21     # optional proxy via httpx (new SDK
22     # does NOT support 'proxies=' kwarg)
23     proxy = os.getenv("HTTPS_PROXY") or os
24     .getenv("HTTP_PROXY")
25     if os.getenv("DOUCHAT_DISABLE_PROXY"
26     ) == "1" or not proxy:
27         return OpenAI(api_key=api_key)
28
29     http_client = httpx.Client(proxies=
30         proxy, timeout=30.0)
```

```
27     return OpenAI(api_key=api_key,
28                     http_client=http_client)
29
30 @router.get("/")
31 def health():
32     status = {
33         "database": "unknown",
34         "chroma": "unknown",
35         "openai": "unknown",
36         "openai_module": getattr(
37             openai_pkg, "__file__", "n/a"),
38         "openai_version": getattr(
39             openai_pkg, "__version__", "n/a"),
40         "collection": os.getenv(
41             "CHROMA_COLLECTION", "docuchat"),
42         "model": os.getenv("OPENAI_MODEL"
43             , "gpt-4o-mini"),
44         "impl": "routers.health.v2", # so
45             we know THIS file is live
46     }
47
48     # 1) DB check
49     try:
50         with engine.connect() as conn:
51             conn.execute(text("SELECT 1"))
52             status["database"] = "ok"
53     except Exception as e:
54         status["database"] = f"error: {e!r}"
55
56     # 2) Chroma check
57     try:
58         coll = rag._ensure_collection()
```

```
53         _ = coll.count()
54         status["chroma"] = "ok"
55     except Exception as e:
56         status["chroma"] = f"error: {e!r}"
57
58     # 3) OpenAI check (new SDK only, NO
59     # proxies kwarg anywhere)
60     try:
61         client = _make_openai_client()
62         resp = client.chat.completions.
63         create(
64             model=os.getenv("OPENAI_MODEL",
65             , "gpt-4o-mini"),
66             messages=[{"role": "user", "content": "Say OK"}],
67             max_tokens=1,
68             temperature=0,
69             )
70         _ = resp.choices[0].message.
71         content
72         status["openai"] = "ok"
73     except Exception as e:
74         status["openai"] = f"error: {e!r}"
75
76     return status
77
```

```
1 # reset_user.py
2 # one-time helper to (re)create a user
   with a properly hashed password
3
4 from docuchat_backend.app.database import
  SessionLocal
5 from docuchat_backend.app import models
6 from docuchat_backend.app.auth import
  get_password_hash
7
8 EMAIL = "umit.ipek@gmail.com"
9 PLAINTEXT_PASSWORD = "12345678" # this is
   what you'll log in with
10
11 db = SessionLocal()
12
13 try:
14     user = db.query(models.User).filter(
15         models.User.email == EMAIL).first()
16     if user:
17         print(f"User {EMAIL} exists.
18             Updating password hash.")
19         user.password_hash =
20             get_password_hash(PLAINTEXT_PASSWORD)
21     else:
22         print(f"User {EMAIL} not found.
23             Creating new user.")
24         user = models.User(
25             email=EMAIL,
26             password_hash=
27                 get_password_hash(PLAINTEXT_PASSWORD),
28             )
29         db.add(user)
```

```
25
26     db.commit()
27     print("Done. You can now log in with
28         that email + password.")
29 finally:
30     db.close()
31
```

```
1 import streamlit as st
2 import requests
3 import json
4 import os
5 from typing import Optional, List
6
7 API_BASE = "http://127.0.0.1:8000"
8
9 st.set_page_config(page_title="DocuChat",
10                     layout="centered")
11
12 # --- session state helpers ---
13 def init_session():
14     if "token" not in st.session_state:
15         st.session_state.token = None
16     if "documents" not in st.session_state:
17         st.session_state.documents = []
18     if "selected_doc_ids" not in st.
19         session_state:
20         st.session_state.selected_doc_ids
21         = []
22     if "health" not in st.session_state:
23         st.session_state.health = None
24
25 def do_login(email: str, password: str
26             ) -> Optional[str]:
27     """
28     Call /auth/login and return the JWT
29     access token if successful, else None.
30     """
31
32     response = requests.post(
33         f"{API_BASE}/auth/login",
34         json={"email": email, "password": password})
35
36     if response.status_code == 200:
37         return response.json().get("token")
38     else:
39         return None
```

```
28     try:
29         resp = requests.post(
30             f"{API_BASE}/auth/login",
31             data={"username": email, "
32                 password": password},
33                 headers={"Content-Type": "
34                     application/x-www-form-urlencoded"}, "
35                     timeout=30,
36                 )
37         if resp.status_code == 200:
38             data = resp.json()
39             return data.get("access_token"
40         )
41         else:
42             st.error(f"Login failed ({resp
43 .status_code}): {resp.text}")
44         return None
45     except Exception as e:
46         st.error(f"Login error: {e}")
47         return None
48
49
50 def fetch_documents(token: str) -> List[
51     dict]:
52     """
53     Fetch the current user's documents
54     from /documents/
55     """
56     try:
57         resp = requests.get(
58             f"{API_BASE}/documents/",
59             headers={"Authorization": f"
60                 Bearer {token}"},
```

```
54                 timeout=30,
55             )
56             if resp.status_code == 200:
57                 return resp.json()
58             else:
59                 st.error(f"Failed to load
60             documents ({resp.status_code}): {resp.text
61             }")
62         except Exception as e:
63             st.error(f"Error fetching
64             documents: {e}")
65         return []
66
67
68 def fetch_health(token: str) -> dict:
69     """
70     Get backend health (/health/) to show
71     model/chroma/db status.
72     """
73     try:
74         resp = requests.get(
75             f"{API_BASE}/health/",
76             headers={"Authorization": f"
77             Bearer {token}"},
78             timeout=10,
79         )
80         if resp.status_code == 200:
81             return resp.json()
82         else:
83             return {"error": f"{resp.
84             status_code} {resp.text}"}
85     except Exception as e:
```

```
81         return {"error": str(e)}
```

```
82
```

```
83
```

```
84 def upload_file(token: str, uploaded_file
85     ) -> Optional[dict]:
86     """
87     Send file -> /documents/upload with
88     long timeout for big PDFs.
89     """
90     headers = {
91         "Authorization": f"Bearer {token}"
92     }
93     files = {
94         "file": (
95             uploaded_file.name,
96             uploaded_file.getvalue(),
97             uploaded_file.type or "
98             application/octet-stream",
99         )
100    }
101   try:
102       resp = requests.post(
103           f"{API_BASE}/documents/upload",
104           headers=headers,
105           files=files,
106           timeout=300, # bumped so 10-
107           # 30MB PDFs won't time out
108       )
109       if resp.status_code == 201:
```

```
108             return resp.json()
109         else:
110             st.error(f"Upload failed ({resp.status_code}): {resp.text}")
111         return None
112     except requests.exceptions.
113         ReadTimeout:
114             st.error("Upload timed out while
115             waiting for the server (still processing
116             ?). Try again or use a smaller file.")
117             return None
118
119
120 def ask_question(token: str, question:
121     str, doc_ids: List[int], top_k: int = 3
122     ) -> Optional[dict]:
123     """
124     Call /query (POST) with question,
125     selected document IDs, top_k.
126     """
127     headers = {
128         "Authorization": f"Bearer {token}"
129     },
130         "Content-Type": "application/json"
131     },
132     }
133
134     body = {
135         "question": question,
136         "document_ids": doc_ids,
```

```
132         "top_k": top_k,
133     }
134
135     try:
136         resp = requests.post(
137             f"{API_BASE}/query/",
138             headers=headers,
139             data=json.dumps(body),
140             timeout=60,
141         )
142         if resp.status_code == 200:
143             return resp.json()
144         else:
145             st.error(f"Query failed ({resp.status_code}): {resp.text}")
146         return None
147     except Exception as e:
148         st.error(f"Query error: {e}")
149     return None
150
151
152 # ----- UI starts here
-----  

153 init_session()
154
155 st.title("DocuChat   ")
156
157 # If not logged in: show login form only
158 if not st.session_state.token:
159     st.subheader("Login")
160
161     email = st.text_input("Email", value=
"umit.ipek@gmail.com", autocomplete="
```

```
161 username")
162     password = st.text_input("Password",
163                             type="password", value="12345678",
164                             autocomplete="current-password")
165
166     if st.button("Sign in"):
167         token = do_login(email, password)
168         if token:
169             st.session_state.token =
170                 token
171             # fetch docs + health once on
172             login
173             st.session_state.documents =
174                 fetch_documents(token)
175             st.session_state.health =
176                 fetch_health(token)
177             st.success("Logged in ✅")
178             st.rerun()
179
180             st.stop()
181
182             # If logged in: main app
183             tab_files, tab_ask, tab_status = st.tabs([
184                 ["📁 Files", "? Ask", ">Status"])
185
186             # --- Tab 1: Files ---
187             with tab_files:
188                 st.header("Upload & Manage Files")
189
190                 uploaded_file = st.file_uploader(
191                     "Choose a file to upload (PDF,
```

```
185 TXT, DOCX). Up to ~30MB.",  
186         type=["pdf", "txt", "docx", "md"  
187     ],  
188     )  
189     if st.button("Upload file"):  
190         if not uploaded_file:  
191             st.warning("Please select a  
file first.")  
192         else:  
193             with st.spinner("Uploading &  
indexing... this can take a minute for  
big PDFs 🚧"):  
194                 result = upload_file(st.  
session_state.token, uploaded_file)  
195                 if result:  
196                     st.success(f"Uploaded: {  
result.get('filename')} (id={result.get('  
id')})")  
197                     # Refresh doc list  
198                     st.session_state.  
documents = fetch_documents(st.  
session_state.token)  
199  
200         st.subheader("Your documents")  
201         if not st.session_state.documents:  
202             st.info("No documents yet.")  
203         else:  
204             # make a lookup dict for display  
205             doc_label_map = {  
206                 d["id"]: f"{d['id']} - {d['  
filename']} ({d['status']})"  
207             for d in st.session_state.
```

```
207 documents
208     }
209
210     # multiselect for which docs to
211     # use in questions
211     current_selection_labels = [
212         doc_label_map.get(doc_id, f"{
213             doc_id}")
213         for doc_id in st.
214             session_state.selected_doc_ids
214         ]
215
216     chosen = st.multiselect(
217         "Select documents to use for
218         Q&A",
218         options=list(doc_label_map.
219             values()),
219         default=
220             current_selection_labels,
220         )
221
222     # reverse map labels -> ids
223     chosen_ids = []
224     for label in chosen:
225         # label format: "18 - file.
225         pdf (ready)"
226             first_part = label.split(
226                 " - ")[0].strip()
227             try:
228                 chosen_ids.append(int(
228                     first_part))
229             except:
230                 pass
```

```
231
232         st.session_state.selected_doc_ids
233         = chosen_ids
234
235         st.write("Currently selected doc
236         IDs:", st.session_state.selected_doc_ids)
237
238 # --- Tab 2: Ask ---
239
240 with tab_ask:
241     st.header("Ask Your Documents")
242
243     if not st.session_state.
244         selected_doc_ids:
245             st.info("Select one or more
246             documents in the Files tab first.")
247     else:
248         question = st.text_area("Your
249         question:", "Summarize in two sentences."
250         )
251
252         top_k = st.slider("How many
253         relevant chunks to retrieve (top_k)", 1,
254         8, 3)
255
256         if st.button("Ask"):
257             with st.spinner("Thinking..."):
258                 answer = ask_question(
259                     st.session_state.
260                     token,
261                     question,
262                     st.session_state.
263                     selected_doc_ids,
264                     top_k=top_k,
```

```
253                     )
254             if answer:
255                 st.subheader("Answer")
256                 st.write(answer.get(
257                     "answer", ""))
258                 st.subheader("Citations")
259                 st.json(answer.get(
260                     "citations", []))
261 # --- Tab 3: Status ---
262 with tab_status:
263     st.header("Backend status")
264     st.write("This checks `/health/` on
265         the FastAPI server.")
266     if st.button("Refresh status"):
267         st.session_state.health =
268             fetch_health(st.session_state.token)
269     if st.session_state.health:
270         st.json(st.session_state.health)
271     else:
272         st.info("No status yet. Click
273             Refresh status.")
```

```
1 fastapi==0.115.0
2 uvicorn==0.32.0
3 sqlalchemy==2.0.36
4 python-jose==3.3.0
5 passlib[bcrypt]==1.7.4
6 chromadb==0.4.24
7 sentence-transformers==3.2.1
8 pdfplumber==0.11.3
9 python-docx==1.1.2
10 openai==1.52.2
11 pydantic==2.9.2
12 python-multipart==0.0.12
13 tqdm==4.67.1
14 anyio==4.11.0
15 typing-extensions>=4.11
16
```

```
1 """DocuChat FastAPI application entry point.
2
3 This module creates the FastAPI app, includes routers and sets up
4 middleware. To run locally, execute ``
5 uvicorn app.main:app --reload``.
6 """
7
8 from fastapi import FastAPI
9
10 from .routers import auth as auth_router
11 from .routers import documents as documents_router
12 from .routers import query as query_router
13 from .routers import health as health_router
14
15 app = FastAPI(title="DocuChat API",
16                 version="0.1.0")
17
18 # Include API routers
19 app.include_router(auth_router.router)
20 app.include_router(documents_router.router)
21 app.include_router(query_router.router)
22 app.include_router(health_router.router)
23
24 @app.get("/", tags=["health"])
25 def read_root():
26     """Health check endpoint."""
27     return {"status": "ok"}  
28
```

```
1 # docuchat_backend/app/auth.py
2 from __future__ import annotations
3
4 import os
5 import logging
6 from datetime import datetime, timedelta,
7     timezone
8 from typing import Optional
9
10 from fastapi import Depends, HTTPException,
11     status
12 from fastapi.security import HTTPBearer,
13     HTTPAuthorizationCredentials
14 from jose import jwt, JWTError
15 from passlib.context import CryptContext
16 from sqlalchemy.orm import Session
17
18 from .database import SessionLocal
19 from . import models
20
21 # -----
22 # Password hashing
23 # -----
24 pwd_context = CryptContext(schemes=["bcrypt"], deprecated="auto")
25
26 def verify_password(plain_password: str,
27     hashed_password: str) -> bool:
28     return pwd_context.verify(
```

```
27 plain_password, hashed_password)
28
29 def get_password_hash(password: str) ->
str:
30     return pwd_context.hash(password)
31
32 # -----
33 # DB dependency
34 # -----
35 def get_db():
36     db = SessionLocal()
37     try:
38         yield db
39     finally:
40         db.close()
41
42 # -----
43 # JWT config
44 # -----
45 SECRET_KEY = os.getenv("JWT_SECRET", "
change-this-in-production")
46 ALGORITHM = os.getenv("JWT_ALG", "HS256")
47 ACCESS_TOKEN_EXPIRE_MINUTES = 60 # 1 hour
48
49 def create_access_token(data: dict,
    expires_delta: Optional[timedelta] = None
) -> str:
50     to_encode = data.copy()
51     expire = datetime.now(timezone.utc
) + (
52         expires_delta if expires_delta
else timedelta(minutes=
ACCESS_TOKEN_EXPIRE_MINUTES)
```

```
53     )
54     to_encode.update({"exp": expire})
55     encoded_jwt = jwt.encode(to_encode,
56                               SECRET_KEY, algorithm=ALGORITHM)
57
58 # -----
59 # Authenticate user for /auth/login
60 # -----
61 def authenticate_user(db: Session, email: str, password: str) -> Optional[models.User]:
62     user = db.query(models.User).filter(
63         models.User.email == email).first()
64     if not user:
65         return None
66     if not verify_password(password, user.password_hash):
67         return None
68
69 # -----
70 # Bearer token dependency for protected
71 # routes
72 bearer_scheme = HTTPBearer(auto_error=False)
73
74 def get_current_user(
75     db: Session = Depends(get_db),
76     credentials:
77     HTTPAuthorizationCredentials = Depends(
78         bearer_scheme),
```

```
77 ) -> models.User:  
78     """  
79         - Pulls Authorization header ("Bearer  
80             <token>")  
81             - Decodes the JWT  
82             - Loads that user from DB  
83             - Raises 401 if anything fails  
84     """  
85     if credentials is None:  
86         logger.warning("get_current_user  
87             : no credentials provided")  
88         raise HTTPException(  
89             status_code=status.  
90                 HTTP_401_UNAUTHORIZED,  
91                     detail="Not authenticated",  
92                     )  
93  
94     try:  
95         payload = jwt.decode(token,  
96             SECRET_KEY, algorithms=[ALGORITHM])  
97         logger.info(f"get_current_user:  
98             decoded payload {payload}")  
99     except JWTError as e:  
100        logger.warning(f"get_current_user  
101            : JWT decode failed: {e}")  
102        raise HTTPException(  
103            status_code=status.  
104                HTTP_401_UNAUTHORIZED,  
105                    detail="Invalid token",
```

```
102 )
103
104     email: str = payload.get("sub")
105     if email is None:
106         logger.warning("get_current_user
107 : 'sub' missing in JWT payload")
108         raise HTTPException(
109             status_code=status.
110             HTTP_401_UNAUTHORIZED,
111             detail="Invalid token payload
",
112             )
113
114     user = db.query(models.User).filter(
115         models.User.email == email).first()
116     if not user:
117         logger.warning(f"get_current_user
: user {email} not found in DB")
118         raise HTTPException(
119             status_code=status.
120             HTTP_401_UNAUTHORIZED,
121             detail="User not found",
122             )
```

```
1 """SQLAlchemy models for DocuChat.
2
3 Defines the database schema for users,
4 documents, collections and chat
5 history. Uses SQLAlchemy's declarative
6 base from ``app.database``.
7 """
8
9 from datetime import datetime
10 from sqlalchemy import Column, Integer,
11     String, ForeignKey, Text, DateTime
12 from sqlalchemy.orm import relationship
13 from .database import Base
14
15 class User(Base):
16     __tablename__ = "users"
17     id = Column(Integer, primary_key=True,
18                 index=True)
19     email = Column(String(255), unique=True,
20                    nullable=False, index=True)
21     password_hash = Column(String(255),
22                            nullable=False)
23     created_at = Column(DateTime, default=
24                           datetime.utcnow)
25     last_login = Column(DateTime, nullable=True)
26     documents = relationship("Document",
27                               back_populates="owner", cascade="all,
28                               delete-orphan")
29     collections = relationship("Collection",
30                                back_populates="owner", cascade="all,
31                                delete-orphan")
32     chat_history = relationship("ChatHistory",
33                                 back_populates="owner",
```

```
20 cascade="all, delete-orphan")
21
22 class Collection(Base):
23     __tablename__ = "collections"
24     id = Column(Integer, primary_key=True
25 , index=True)
25     user_id = Column(Integer, ForeignKey(
26         "users.id", ondelete="CASCADE"), nullable=
27 False)
26     name = Column(String(255), nullable=
28 False)
27     created_at = Column(DateTime, default=
28         datetime.utcnow)
28     owner = relationship("User",
29         back_populates="collections")
29     documents = relationship("Document",
30         back_populates="collection", cascade="all
31 , delete-orphan")
30
31 class Document(Base):
32     __tablename__ = "documents"
33     id = Column(Integer, primary_key=True
34 , index=True)
34     user_id = Column(Integer, ForeignKey(
35         "users.id", ondelete="CASCADE"), nullable=
36 False)
35     filename = Column(String(255),
36         nullable=False)
36     file_path = Column(String(500),
37         nullable=False)
37     file_size = Column(Integer, nullable=
38 True)
38     file_type = Column(String(50),
```

```
38 nullable=True)
39     status = Column(String(50), default="processing")
40     summary = Column(Text, nullable=True)
41     upload_date = Column(DateTime, default=datetim
e.utcnow)
42     collection_id = Column(Integer,
43     ForeignKey("collections.id", ondelete="SET
44     NULL"), nullable=True)
45     owner = relationship("User",
46     back_populates="documents")
47     collection = relationship("Collection"
48     , back_populates="documents")
49     chat_history = relationship(
50     "ChatHistory", back_populates="document",
51     cascade="all, delete-orphan")
52
53 class ChatHistory(Base):
54     __tablename__ = "chat_history"
55     id = Column(Integer, primary_key=True
56     , index=True)
57     user_id = Column(Integer, ForeignKey(
58     "users.id", ondelete="CASCADE"), nullable=
59     False)
60     document_id = Column(Integer,
61     ForeignKey("documents.id", ondelete="CASCAD
E"), nullable=False)
62     question = Column(Text, nullable=False
63     )
64     answer = Column(Text, nullable=False)
65     sources = Column(Text, nullable=True)
66     rating = Column(Integer, nullable=True
67     )
```

```
56     created_at = Column(DateTime, default=
      datetime.utcnow)
57     owner = relationship("User",
      back_populates="chat_history")
58     document = relationship("Document",
      back_populates="chat_history")
59
```

```
1 """Core configuration and utility  
functions for DocuChat."""  
2
```

```
1 # docuchat_backend/app/routers/query.py
2 from __future__ import annotations
3
4 from fastapi import APIRouter, Depends,
5     HTTPException
6 from pydantic import BaseModel, Field
7 from typing import List, Optional
8
9 from .. import rag, auth, models
10
11 router = APIRouter(prefix="/query", tags=[
12     "query"])
13
14 class QueryRequest(BaseModel):
15     question: str = Field(..., description=
16         "Natural language question to ask your
17         docs")
18     document_ids: Optional[List[int]] =
19         Field(None, description="List of document
20         IDs to restrict search")
21     top_k: int = Field(4, description="How
22         many chunks to retrieve internally")
23
24 class QueryResponse(BaseModel):
25     answer: str
26     citations: Optional[List[dict]]
27
28
29 @router.post("/", response_model=
30     QueryResponse)
31
32 def ask_question(
33     req: QueryRequest,
34     current_user: models.User = Depends(
35         auth.get_current_user),
```

```
25  ):
26      """
27      Auth required. Runs semantic search +
28      LLM answer (if API key is set).
29      """
30      try:
31          result = rag.query_documents(
32              question=req.question,
33              document_ids=req.document_ids,
34              top_k=req.top_k,
35          )
36      except Exception as e:
37          raise HTTPException(status_code=
38              500, detail=str(e))
39      # result should be {"answer": ..., "citations": ...}
40      return QueryResponse(
41          answer=result["answer"],
42          citations=result["citations"],
43      )
```

```
1 """
2 DocuChat Backend Package
3
4 This package contains the FastAPI
5 application and supporting modules for
6 user authentication, database models and
7 API routes.
8 from .main import app # noqa: F401
9
```

```
1 """Database configuration and session
2 management for DocuChat.
3
4 This module defines the SQLAlchemy engine
5 and session maker. It reads
6 configuration from environment variables
7 with sensible defaults, and
8 exports a dependency that yields a
9 database session per request.
10
11 def get_database_url() -> str:
12     """Construct the database URL from
13     environment variables.
14
15     For local development, fall back to an
16     SQLite database so that the
17     backend can start without external
18     dependencies. In production,
19     override `DATABASE_URL` with a
20     PostgreSQL connection string.
21     """
22
23     return os.getenv("DATABASE_URL", "
24         sqlite:///./docuchat.db")
25
26 # Create the SQLAlchemy engine. For
27 # SQLite, we need check_same_thread=False
28 # because FastAPI runs each request in a
29 # separate thread.
```

```
22 DATABASE_URL = get_database_url()
23 connect_args = {} # extra arguments for
create_engine
24 if DATABASE_URL.startswith("sqlite"): #
    pragma: no cover
25     connect_args = {"check_same_thread":  
26         False}
27 engine = create_engine(DATABASE_URL,
28 connect_args=connect_args)
29
30 class Base(DeclarativeBase):
31     """Base class for SQLAlchemy models
32 ."""
33
34 # Create a configured "Session" class
35 SessionLocal = sessionmaker(autocommit=False, autoflush=False, bind=engine)
36
37 def get_db():
38     """FastAPI dependency that provides a
39 database session per request.
40
41     Yields a SQLAlchemy session and
42 ensures it is closed after the request
43     is complete, even if an exception
44 occurs.
45     """
46     db = SessionLocal()
47     try:
48         yield db
49     finally:
50         db.close()
```

```
1 """API routers for DocuChat.  
2  
3 Routers are grouped by feature area. They  
4 are included in the main  
5 application in ``app.main``.  
6 """
```

```
1 """File storage utilities.
2
3 This module abstracts saving uploaded
4 files to the local filesystem. In
5 production you might replace this with an
6 S3 or Google Cloud Storage
7 client. Filenames are randomised to
8 prevent collisions and to avoid
9 disclosing the original filenames on disk.
10 """
11
12 UPLOAD_DIR = os.getenv("UPLOAD_DIR", "uploaded_files")
13
14 def ensure_upload_dir() -> None:
15     """Ensure that the upload directory
16     exists."""
17     os.makedirs(UPLOAD_DIR, exist_ok=True)
18
19 def save_upload_file(upload_file:
20     UploadFile) -> str:
21     """Save an uploaded file and return
22     the file path."""
23     ensure_upload_dir()
24     _, ext = os.path.splitext(upload_file.
25     filename)
26     unique_name = f"{uuid.uuid4().hex}{ext}"
27
28     file_path = os.path.join(UPLOAD_DIR,
29     unique_name)
```

```
24     with open(file_path, "wb") as buffer:
25         while True:
26             contents = upload_file.file.
27             read(1024 * 1024)
28             if not contents:
29                 break
30             buffer.write(contents)
31
32 def delete_file(path: str) -> None:
33     """Delete a file from the filesystem
34     .
35     """
36     try:
37         os.remove(path)
38     except FileNotFoundError:
39         pass
```

```
1 """Application configuration settings.  
2  
3 Values can be overridden via environment  
4 variables. For example,  
5 ``SECRET_KEY`` controls JWT token signing  
6 , and ``ACCESS_TOKEN_EXPIRE_MINUTES``  
7 defines how long login tokens remain valid  
8 .  
9 """  
10 import os  
11  
12 class Settings:  
13     # Cryptographic key used to sign JWT  
14     # tokens. In production, set this to a  
15     # long, random string via an  
16     # environment variable.  
17     SECRET_KEY: str = os.getenv("  
18         SECRET_KEY", "CHANGEME_SUPER_SECRET")  
19     ALGORITHM: str = os.getenv("  
20         JWT_ALGORITHM", "HS256")  
21     ACCESS_TOKEN_EXPIRE_MINUTES: int = int  
22     (os.getenv("ACCESS_TOKEN_EXPIRE_MINUTES",  
23         "30"))  
24  
25 settings = Settings()  
26
```

```
1 # docuchat_backend/app/rag.py
2 from __future__ import annotations
3
4 import os
5 import io
6 import logging
7 from typing import List, Optional, Dict,
8 Any
9
10 logger = logging.getLogger("docuchat.rag")
11 logger.setLevel(logging.INFO)
12
13 # ----- external libs (optional
14 # imports with graceful fallback) -----
15
16 # text PDFs
17 try:
18     import pdfplumber
19 except Exception:
20     pdfplumber = None # type: ignore
21
22 # docx
23 try:
24     import docx # python-docx
25 except Exception:
26     docx = None # type: ignore
27
28 # PyPDF2 as a second text extractor
29 try:
30     import PyPDF2
31 except Exception:
32     PyPDF2 = None # type: ignore
```

```
32 # OCR stack
33 try:
34     import pytesseract
35 except Exception:
36     pytesseract = None # type: ignore
37
38 try:
39     from pdf2image import
        convert_from_path
40 except Exception:
41     convert_from_path = None # type:
        ignore
42
43 try:
44     from PIL import Image # pillow
45 except Exception:
46     Image = None # type: ignore
47
48 # chroma
49 try:
50     import chromadb
51 except Exception as e:
52     chromadb = None # type: ignore
53     logger.error("Chroma import failed: %s
", e)
54
55 # embeddings
56 try:
57     from sentence_transformers import
        SentenceTransformer
58 except Exception as e:
59     SentenceTransformer = None # type:
        ignore
```

```
60     logger.error("sentence-transformers  
       import failed: %s", e)  
61  
62 # OpenAI chat model  
63 _OPENAI_AVAILABLE = False  
64 try:  
65     from openai import OpenAI  
66     _OPENAI_AVAILABLE = True  
67 except Exception:  
68     _OPENAI_AVAILABLE = False  
69  
70  
71 # ----- env / config -----  
72  
73 CHROMA_DB_DIR = os.getenv("CHROMA_DB_DIR"  
    , "./docuchat_backend/chroma_db")  
74 CHROMA_COLLECTION = os.getenv("CHROMA_COLLECTION", "docuchat")  
75 EMBEDDING_MODEL_NAME = os.getenv("EMBEDDING_MODEL", "sentence-transformers/  
all-MiniLM-L6-v2")  
76  
77 OPENAI_API_KEY = os.getenv("OPENAI_API_KEY")  
78 OPENAI_MODEL = os.getenv("OPENAI_MODEL", "gpt-4o-mini")  
79  
80 # IMPORTANT: update these two for your  
# Windows install paths.  
81 # If you add tesseract/poppler to PATH  
# later, you can set them to None.  
82 TESSERACT_EXE = r"C:\Program Files\  
Tesseract-OCR\tesseract.exe"
```

```
83 POPLER_BIN_DIR = r"C:\Tools\poppler-25.07
.0\Library\bin" # <- example; update to
your actual poppler bin folder
84
85 # ----- globals -----
86 _COLLECTION = None
87 _S_BERT = None
88
89
90 # ----- helpers: embeddings
-----
91
92 def _load_sbert() -> SentenceTransformer
| None:
93     """lazy-load embedding model"""
94     global _S_BERT
95     if _S_BERT is not None:
96         return _S_BERT
97     if SentenceTransformer is None:
98         logger.warning("sentence-
transformers not available; RAG will
degrade.")
99     return None
100    logger.info(f"Loading embeddings
model: {EMBEDDING_MODEL_NAME}")
101    _S_BERT = SentenceTransformer(
102        EMBEDDING_MODEL_NAME)
103    return _S_BERT
104
105 def _embed_texts(texts: List[str]) ->
List[List[float]]:
106    """
```

```
107      Turn text chunks into embeddings via
108          sentence-transformers.
109      """
110      model = _load_sbert()
111      if model is None:
112          logger.warning("Embeddings model
113              missing; using zero vectors.")
114      return [[0.0] * 384 for _ in
115          texts]
116
117      vecs = model.encode(
118          texts,
119          convert_to_tensor=False,
120          normalize_embeddings=True,
121          )
122
123 # ----- helpers: chroma -----
124
125 def _ensure_collection():
126     """Return / create persistent Chroma
127         collection."""
128     if chromadb is None:
129         raise RuntimeError("ChromaDB is
130             not installed.")
131
132     global _COLLECTION
133     if _COLLECTION is not None:
134         return _COLLECTION
```

```
134     os.makedirs(CHROMA_DB_DIR, exist_ok=True)
135     logger.info(f"Opening Chroma at {CHROMA_DB_DIR!r}, collection={CHROMA_COLLECTION!r}")
136     client = chromadb.PersistentClient(path=CHROMA_DB_DIR)
137
138     try:
139         collection = client.get_collection(CHROMA_COLLECTION)
140     except ValueError:
141         collection = client.create_collection(CHROMA_COLLECTION)
142
143     _COLLECTION = collection
144     return _COLLECTION
145
146
147 # ----- helpers: text splitting
148 # -----
149 def _split_text(text: str, chunk_size: int = 900, chunk_overlap: int = 150) -> List[str]:
150     text = text.strip().replace("\r\n", "\n")
151     if not text:
152         return []
153     out: List[str] = []
154     start = 0
155     n = len(text)
156     while start < n:
```

```
157         end = min(n, start + chunk_size)
158         chunk = text[start:end]
159         out.append(chunk)
160         if end == n:
161             break
162         start = max(end - chunk_overlap,
163                     start + 1)
164
165
166 # ----- helpers: PDF extraction
167 # paths -----
168 def _extract_pdf_text_pdfplumber(path: str) -> str:
169     """Try pdfplumber to pull machine-
170     readable text."""
171     if pdfplumber is None:
172         return ""
173     try:
174         parts: List[str] = []
175         with pdfplumber.open(path) as pdf
176             :
177                 for page in pdf.pages:
178                     try:
179                         t = page.extract_text()
180                     except Exception:
181                         t = ""
182                     if t.strip():
183                         parts.append(t)
184     return "\n\n".join(parts).strip()
185 except Exception as e:
```

```
184         logger.warning("pdfplumber failed  
185             on %s: %s", path, e)  
186             return ""  
187  
188 def _extract_pdf_text_pypdf2(path: str  
189     ) -> str:  
190     """Fallback using PyPDF2 .  
191     extract_text(). Good for some digital  
192     PDFs."""  
193     if PyPDF2 is None:  
194         return ""  
195     try:  
196         parts: List[str] = []  
197         with open(path, "rb") as f:  
198             reader = PyPDF2.PdfReader(f)  
199             for page in reader.pages:  
200                 try:  
201                     t = page.extract_text  
202                     () or ""  
203                     except Exception:  
204                         t = ""  
205                         if t.strip():  
206                             parts.append(t)  
207                         return "\n\n".join(parts).strip()  
208             except Exception as e:  
209                 logger.warning("PyPDF2 failed on  
210                 %s: %s", path, e)  
211                 return ""  
212  
213 def _ocr_pdf_pages(path: str) -> str:  
214     """
```

```
211      LAST RESORT: OCR every page image
212      with Tesseract.
213      Requires:
214          - pdf2image.convert_from_path()
215          - pytesseract
216          - Tesseract binary
217          - Poppler (for pdf2image on Windows)
218          """
219      if convert_from_path is None or
220          pytesseract is None:
221          return ""
222      # configure external binaries if
223      # needed
224      # (if you added them to PATH already
225      # , you can skip this)
226      if TESSERACT_EXE and hasattr(
227          pytesseract, "pytesseract"):
228          pytesseract.pytesseract.
229          tesseract_cmd = TESSERACT_EXE # point to
230          tesseract.exe
231          poppler_kw = {}
232          if POPLER_BIN_DIR and
233              convert_from_path is not None:
234              poppler_kw["poppler_path"] =
235                  POPLER_BIN_DIR
236
237          try:
238              # convert PDF pages -> list of
239              # PIL Images
240              images = convert_from_path(path
```

```
File - C:\Users\uipek\PycharmProjects\DocuChat\docuchat_backend\app\rag.py
233 , **poppler_kw) # this can be slow for
    huge PDFs
234     except Exception as e:
235         logger.warning("pdf2image failed
    to rasterize %s: %s", path, e)
236     return ""
237
238     ocr_parts: List[str] = []
239     for i, img in enumerate(images):
240         try:
241             text = pytesseract.
    image_to_string(img, lang="eng") or ""
242         except Exception as e:
243             logger.warning("pytesseract
    failed on page %s of %s: %s", i, path, e)
244         text = ""
245         if text.strip():
246             ocr_parts.append(text.strip()
    ())
247
248     return "\n\n".join(ocr_parts).strip()
249
250
251 # ----- helpers: DOCX / fallback
    -----
252
253 def _read_docx(path: str) -> str:
254     if docx is None:
255         raise RuntimeError("python-docx
    not installed; run: pip install python-
    docx")
256     try:
257         d = docx.Document(path)
```

File - C:\Users\uipek\PycharmProjects\DocuChat\docuchat\_backend\app\rag.py

```
258         txt = "\n".join(p.text for p in d
259                         .paragraphs)
260     return txt.strip() or "[  
261         NO_EXTRACTED_TEXT_FROM_DOCX]"
262     except Exception as e:
263         logger.warning("DOCX read failed  
264             for %s: %s", path, e)
265     return "[DOCX_READ_ERROR]"
266
267 # ----- master extractor -----
268
269 def _read_file_text(file_path: str) ->
270     str:
271     """
272     Unified extractor with OCR fallback.  

273     1. direct text (.txt, .md, .log)  

274     2. pdfplumber -> PyPDF2 -> OCR  

275     3. docx  

276     4. best-effort bytes decode
277     """
278     ext = os.path.splitext(file_path)[1].lower()
279
280     # plain text-ish
281     if ext in [".txt", ".md", ".log"]:
282         with open(file_path, "r",
283                    encoding="utf-8", errors="ignore") as f:
284             return f.read().strip()
285
286     # pdf
287     if ext == ".pdf":
288         # step 1: machine-readable
```

```
285         text =
286         _extract_pdf_text_pdfplumber(file_path)
287         if not text.strip():
288             text =
289             _extract_pdf_text_pypdf2(file_path)
290
291             # step 2: OCR fallback
292             if (not text.strip()) or text.
293                 strip().startswith("[NO_EXTRACTED_TEXT"):
294                     logger.info("Attempting OCR
295                     for PDF with little/no text: %s",
296                     file_path)
297                     ocr_text = _ocr_pdf_pages(
298                     file_path)
299                     if ocr_text.strip():
300                         text = ocr_text
301
302                     if not text.strip():
303                         text = "[
304                             NO_EXTRACTED_TEXT_FROM_PDF]"
305
306                     return text
307
308                     # word
309                     if ext in [".docx", ".doc"]:
310                         return _read_docx(file_path)
311
312
313                     # fallback raw bytes
314                     with open(file_path, "rb") as f:
315                         raw = f.read()
316                         for enc in ("utf-8", "latin-1"):
317                             try:
318                                 decoded = raw.decode(enc)
319                                 if decoded.strip():


```

```
311             return decoded
312     except Exception:
313         pass
314
315     return "[BINARY_OR_EMPTY]"
316
317
318 # ----- answer helpers -----
319
320 def _build_prompt(question: str, contexts
321 : List[str]) -> str:
322     ctx = "\n\n---\n\n".join(contexts)
323     return (
324         "You are a helpful assistant.
325         Answer ONLY from the provided context. "
326         "If the answer is not in the
327         context, say you don't know.\n\n"
328         f"Context:{ctx}\n\n"
329         f"Question: {question}\n\n"
330         "Answer:"
331     )
332
333
334
335     client = OpenAI(api_key=
336         OPENAI_API_KEY)
337     resp = client.chat.completions.create
338     (
339         model=OPENAI_MODEL,
```

```
337         messages=[{"role": "user", "content": prompt}],
338         temperature=0.2,
339     )
340     return resp.choices[0].message.
341     content.strip()
342
343 def _answer_extractive(hits: List[Dict[
344     str, Any]]) -> str:
345     parts = []
346     for h in hits:
347         meta = h.get("metadata", {})
348         idx = meta.get("chunk_index", "?")
349
350         txt = h.get("document", "").strip()
351         if len(txt) > 600:
352             txt = txt[:600].rstrip() +
353             """
354             parts.append(f"\n{txt}\n[cite: doc={did} chunk={idx}]")
355
356     if not parts:
357         return "I couldn't find relevant
358         context in the selected documents."
359
360     return "\n\n---\n\n".join(parts)
361
```

File - C:\Users\uipek\PycharmProjects\DocuChat\docuchat\_backend\app\rag.py

```
362 # ----- public API -----
363
364 def store_document(doc_id: int, file_path
365     : str) -> None:
366     """
367         Read file, split text, embed, upsert
368         into Chroma.
369         Always index something (even OCR text
370         or placeholder) so status can be 'ready'
371         .
372     """
373     collection = _ensure_collection()
374
375     full_text = _read_file_text(file_path
376 )
377
378     chunks = _split_text(full_text)
379     if not chunks:
380         chunks = [full_text]
381
382     embeddings = _embed_texts(chunks)
383
384     ids = [f"doc:{doc_id}:chunk:{i}" for
385         i in range(len(chunks))]
386     metadata = [
387         {
388             "doc_id": doc_id,
389             "chunk_index": i,
390             "filename": os.path.basename(
391                 file_path),
392         }
393         for i in range(len(chunks))
394     ]
395
396     collection.upsert(
397         [{"id": id, "text": chunk, "embedding": embedding}
398          for id, chunk, embedding in zip(ids, chunks, embeddings)])
399
400     return None
```

```
388
389     logger.info(f"Upserting {len(chunks)}"
390                 f"chunks for doc_id={doc_id}")
391     collection.upsert(
392         ids=ids,
393         documents=chunks,
394         metadatas=metadatas,
395         embeddings=embeddings,
396     )
397
398 def query_documents(question: str,
399                      document_ids: Optional[List[int]] = None
400                      , top_k: int = 4) -> Dict[str, Any]:
401     """
402         Retrieve relevant chunks from Chroma
403         , then answer with OpenAI if available,
404         else fall back to stitched extractive
405         answer.
406     """
407     collection = _ensure_collection()
408
409     where: Dict[str, Any] = {}
410
411     if document_ids:
412         where = {"doc_id": {"$in": document_ids}}
413
414     q_emb = _embed_texts([question])[0]
415
416     res = collection.query(
417         query_embeddings=[q_emb],
418         n_results=max(1, top_k),
419         where=where or None,
```

```
415         include=["documents", "metadata"
416             , "distances"],
417         )
418     documents: List[str] = res.get("documents", [[]])[0] if res.get("documents") else []
419     metadata: List[Dict[str, Any]] = res.get("metadata", [[]])[0] if res.get("metadata") else []
420     distances: List[float] = res.get("distances", [[]])[0] if res.get("distances") else []
421
422     hits = []
423     for doc, meta, dist in zip(documents, metadata, distances):
424         hits.append({"document": doc, "metadata": meta, "distance": float(dist)})
425
426     contexts = [h["document"] for h in hits]
427
428     if OPENAI_API_KEY and
429         _OPENAI_AVAILABLE and contexts:
430         try:
431             answer_text =
432                 _answer_with_openai(question, contexts)
433         except Exception as e:
434             logger.warning("OpenAI call failed (%s); falling back.", e)
435             answer_text =
```

```
433     _answer_extractive(hits)
434     else:
435         answer_text = _answer_extractive(
436             hits)
437         citations = [
438             {
439                 "doc_id": h["metadata"].get("doc_id"),
440                 "chunk_index": h["metadata"].get("chunk_index"),
441                 "filename": h["metadata"].get("filename"),
442                 "distance": h["distance"],
443             }
444             for h in hits
445         ]
446
447     return {"answer": answer_text,
448             "citations": citations}
```

```
1 """Authentication endpoints.
2
3 Provides routes for user registration and
4 login, and a dependency for
5 retrieving the currently authenticated
6 user from a JWT token.
7 """
8 from datetime import timedelta
9 from fastapi import APIRouter, Depends,
10    HTTPException, status
11 from fastapi.security import
12    OAuth2PasswordBearer,
13    OAuth2PasswordRequestForm
14 from sqlalchemy.orm import Session
15 from .. import models, schemas, auth
16 from ..database import get_db
17 from ..core.config import settings
18
19 router = APIRouter(prefix="/auth", tags=["auth"])
20
21 # OAuth2PasswordBearer expects token in
22 # Authorization header as: Bearer <token>
23 oauth2_scheme = OAuth2PasswordBearer(
24     tokenUrl="/auth/login")
25
26 @router.post("/register", response_model=
27     schemas.UserRead, status_code=status.
28     HTTP_201_CREATED)
29 def register(user_in: schemas.UserCreate,
30     db: Session = Depends(get_db)):
31     existing = db.query(models.User).
32         filter(models.User.email == user_in.email
33         ).first()
```

```
21     if existing:
22         raise HTTPException(status_code=
23             status.HTTP_400_BAD_REQUEST, detail="Email
24             already registered")
25     hashed_password = auth.
26     get_password_hash(user_in.password)
27     user = models.User(email=user_in.email
28 , password_hash=hashed_password)
29     db.add(user)
30     db.commit()
31     db.refresh(user)
32     return user
33
34 @router.post("/login", response_model=
35     schemas.Token)
36 def login(form_data:
37     OAuth2PasswordRequestForm = Depends(), db
38     : Session = Depends(get_db)):
39     """Authenticate a user and return a
40     JWT token."""
41     user = auth.authenticate_user(db,
42     email=form_data.username, password=
43     form_data.password)
44     if not user:
45         raise HTTPException(status_code=
46             status.HTTP_401_UNAUTHORIZED, detail="
47             Invalid email or password")
48     access_token_expires = timedelta(
49     minutes=settings.
50     ACCESS_TOKEN_EXPIRE_MINUTES)
51     access_token = auth.
52     create_access_token(data={"sub": user.
53     email}, expires_delta=access_token_expires
```

```
37 )
38     return schemas.Token(access_token=
39         access_token, token_type="bearer")
40 def get_current_user(db: Session = Depends(
41     get_db), token: str = Depends(
42     oauth2_scheme)) -> models.User:
43     user = auth.get_user_from_token(db,
44     token)
45     if not user:
46         raise HTTPException(status_code=
47             status.HTTP_401_UNAUTHORIZED, detail="
48             Invalid or expired token")
49     return user
50 
```

```
1 """Background job definitions.  
2  
3 In a production deployment you would use  
4 Celery, RQ or another task  
5 queue to run these functions  
6 asynchronously. For demonstration  
7 purposes, they are synchronous functions  
8 called directly from the API  
9 handlers.  
10  
11 logger = logging.getLogger(__name__)  
12  
13  
14 def process_document_job(db: Session,  
15     document: models.Document) -> None:  
16     """Process a document: extract text,  
17     generate embeddings and update status."""  
18     try:  
19         rag.process_document(doc_id=  
20             document.id, file_path=document.file_path)  
21         document.status = "ready"  
22         db.add(document)  
23         db.commit()  
24     except Exception as exc:  
25         logger.error("Processing failed  
26         for document %s: %s", document.id, exc)  
27         document.status = "error"  
28         db.add(document)  
29         db.commit()
```

```
1 from __future__ import annotations
2
3 from datetime import datetime
4 from typing import List, Optional
5
6 from pydantic import BaseModel, EmailStr,
7     Field
8
9 # =====
10 # Auth / Tokens
11 # =====
12 class Token(BaseModel):
13     access_token: str
14     token_type: str = "bearer"
15     # Optional: seconds until expiry if
16     # you choose to return it
17     expires_in: Optional[int] = None
18
19 # =====
20 # Users
21 # =====
22 class UserBase(BaseModel):
23     email: EmailStr
24
25
26 class UserCreate(UserBase):
27     password: str
28
29
30 class UserRead(UserBase):
31     id: int
```

```
32
33     class Config:
34         # Pydantic v2 replacement for
35         orm_mode
36
37
38 # =====
39 # Documents
40 # =====
41 class DocumentOut(BaseModel):
42     id: int
43     filename: str
44     file_size: Optional[int] = None
45     file_type: Optional[str] = None
46     status: str
47     summary: Optional[str] = None
48     upload_date: datetime
49     collection_id: Optional[int] = None
50
51     class Config:
52         from_attributes = True
53
54
55 # =====
56 # Query (RAG)
57 # =====
58 class QueryRequest(BaseModel):
59     question: str
60     document_ids: Optional[List[int]] =
None
61     # how many chunks to retrieve from the
vector store
```

```
62     top_k: int = Field(default=4, ge=1, le
=20)
63
64
65 class Citation(BaseModel):
66     doc_id: int
67     chunk_index: int
68     filename: Optional[str] = None
69     distance: Optional[float] = None
70
71
72 class QueryResponse(BaseModel):
73     answer: str
74     citations: Optional[List[Citation]] =
None
75
```

```
1 from __future__ import annotations
2
3 import logging
4 from pathlib import Path
5 from typing import List
6
7 from fastapi import APIRouter, Depends,
8     File, HTTPException, UploadFile, status
9 from sqlalchemy.orm import Session
10
11 from .. import models, schemas, auth, rag
12 from ..database import get_db
13
14 logger = logging.getLogger("docuchat.
15 documents")
16
17
18 @router.get("/", response_model=List[
19     schemas.DocumentOut])
20
21 def list_documents(
22     db: Session = Depends(get_db),
23     current_user: models.User = Depends(
24         auth.get_current_user),
25 ):
26     """
27     Return all documents for the logged-in
28     user.
29     """
30     docs = (
31         db.query(models.Document)
```

```
28         .filter(models.Document.user_id
29             == current_user.id)
30         .order_by(models.Document.
31             upload_date.desc())
32     .all()
33
34
35 @router.post("/upload", response_model=
36     schemas.DocumentOut, status_code=status.
37     HTTP_201_CREATED)
38 def upload_document(
39     file: UploadFile = File(...),
40     db: Session = Depends(get_db),
41     current_user: models.User = Depends(
42         auth.get_current_user),
43 ):
44     """
45     Stream the uploaded file to disk,
46     enforce size limit (30MB),
47     insert DB row, then run RAG indexing.
48     """
49     MAX_BYTES = 30 * 1024 * 1024 # 30 MB
50
51     upload_dir = Path(__file__).parent / "uploads"
52     upload_dir.mkdir(parents=True,
53     exist_ok=True)
54
55     safe_name = Path(file.filename).name
56     # strip any path tricks
57     saved_path = upload_dir / safe_name
```

```
52
53     file_size = 0
54     try:
55         with open(saved_path, "wb") as out:
56             while True:
57                 chunk = file.file.read(
58                     1024 * 1024) # 1MB chunk
59                 if not chunk:
60                     break
61                 file_size += len(chunk)
62                 if file_size > MAX_BYTES:
63                     out.close()
64                     saved_path.unlink(
65                         missing_ok=True)
66                     raise HTTPException(
67                         status_code=413,
68                         detail=f"File too
69                         large (> {MAX_BYTES} bytes)",
70                         )
71                     out.write(chunk)
72             finally:
73                 # make sure file cursor reset just
74                 # in case (not strictly needed now)
75                 try:
76                     file.file.seek(0)
77                 except Exception:
78                     pass
79
80             # Create DB row with status=processing
81             doc = models.Document(
82                 user_id=current_user.id,
83                 filename=safe_name,
```

```
80         file_path=str(saved_path),
81         file_size=file_size,
82         file_type=file.content_type or "
83             application/octet-stream",
84         status="processing",
85     )
86     db.add(doc)
87     db.commit()
88     db.refresh(doc)
89
90     # Now try to process & embed
91     try:
92         rag.store_document(doc.id, str(
93             saved_path))
94         doc.status = "ready"
95     except Exception as e:
96         logger.exception("Processing
97 failed for document %s: %s", doc.id, e)
98         doc.status = "error"
99
100    return doc
101
```