

DocuMind AI

Technical Documentation

A Retrieval-Augmented Generation System
for Personal Document Intelligence

Author: Mouadi Yassine

Academic Year: 2025–2026

Computer Science Project
University of The National School of Applied Sciences, Kenitra

Contents

1	Introduction	4
1.1	The Problem with Document Search	4
1.2	Why Retrieval-Augmented Generation (RAG) Works	4
2	Technology Stack	4
2.1	Frontend: React + TypeScript	4
2.2	Backend: FastAPI	5
2.3	Database: PostgreSQL + PGVector	5
2.4	AI: Sentence Transformers + Google Gemini	5
3	Frontend Architecture	5
3.1	Folder Structure	5
3.2	State Management: Lifted State in App.tsx	6
3.3	UI Behavior During Async Operations	7
3.4	File Upload Mechanism	7
4	Backend Architecture	7
4.1	Entry Point: main.py	7
4.2	Router Separation:	8
4.3	Dependency Injection: get_db()	8
5	Authentication System	9
5.1	auth.py as the Receptionist	9
5.2	Password Hashing with bcrypt	9
5.3	Session Management	10
6	Database Design	10
6.1	The Database Defense	10
6.2	Entities Of Our Database:	10
6.3	Foreign Keys with ON DELETE CASCADE	10

6.4	Data Isolation Per User	11
7	AI & Document Processing Pipeline	11
7.1	From PDF Bytes to Text	11
7.2	Chunking Logic: The 500-Character Window	12
7.3	Embedding Generation: Why 384 Dimensions	12
7.4	Singleton Pattern for Model Loading	12
7.5	Semantic Search with Cosine Distance	13
7.6	Question Answering with Gemini	13
8	API Routes in Detail	14
8.1	POST /documents/upload	14
8.2	GET /documents/list/{user_id}	14
8.3	DELETE /documents/{document_id}	14
8.4	DELETE /documents/all/{user_id}	15
8.5	POST /documents/query	15
9	Infrastructure & Deployment	15
9.1	Docker Compose Orchestration	15
9.2	Startup Race Conditions	16
9.3	Database Connection Retry Logic	16
9.4	Frontend Docker Optimizations	16
10	End-to-End Flow: Upload to Query	16
10.1	Step 1: User Selects File	16
10.2	Step 2: Frontend Uploads	17
10.3	Step 3: Backend Receives	17
10.4	Step 4: Read PDF	17
10.5	Step 5: Chunking	17
10.6	Step 6: Database Insert	17
10.7	Step 7: User Asks Question	17

10.8 Step 8: Semantic Search	18
10.9 Step 9: Answer Generation	18
10.10 Step 10: Display	18
11 Error Handling & Resilience	18
11.1 Frontend:	18
11.2 Backend: Structured Errors	18
11.3 Database: Connection Pooling	19
12 Project Structure Recap	19
12.1 Root Infrastructure	19
12.2 Backend Architecture (FastAPI)	19
12.3 Frontend Architecture (React & TypeScript)	20
13 What Works Well	20
13.1 The RAG Pipeline	20
13.2 Docker Compose Setup	20
13.3 Type Safety	20
13.4 Database Design	20
14 Intentional Simplifications	21
14.1 No User Sessions Database	21
14.2 No File Size Limits	21
14.3 No Rate Limiting	21
14.4 No Background Processing	21
15 Conclusion	21

1. Introduction

1.1. The Problem with Document Search

Traditional information retrieval systems are often limited by the constraints of keyword matching, which requires users to recall precise terminology to locate relevant data. While shortcuts like "Ctrl+F" are effective for string searches, they lack the ability to synthesize information or bridge the gap between divergent vocabularies. Consequently, navigating dense academic or professional documentation becomes a fragmented process that relies more on linguistic luck than substantive inquiry.

To address this, semantic search prioritizes intent and context over literal syntax. For example, a student might ask, "What are the consequences of missing the final paper deadline?" while the text may only include "late submission penalties", "integrity policies", etc.

By understanding the conceptual relationship between these terms, semantic search enables users to query complex datasets naturally, ensuring that critical insights are accessible even when the specific terminology is unfamiliar.

1.2. Why Retrieval-Augmented Generation (RAG) Works

Retrieval-Augmented Generation (RAG) offers a more elegant yet simple alternative way to LLMs. Instead of building a giant AI that memorizes everything, which has the obvious inconveniences of being slow and expensive, RAG enables an AI to consult external, verified data sources before generating a response. This methodology functions much like a librarian who:

- You ask a question
- The system searches through your documents for relevant sections
- It sends those sections to the AI along with your question
- The AI answers based only on what it found

This means the AI can't make things up, it cites your actual documents, and it doesn't need retraining when you add new files.

2. Technology Stack

2.1. Frontend: React + TypeScript

Why React? Because components within it are like LEGO bricks. Each piece - Navbar, DocumentUpload, ChatInterface - is isolated, testable, and reusable. When something

breaks, we know exactly which brick to fix. The virtual Document Object Model (DOM) means the UI stays fast even when documents load.

Why TypeScript? TypeScript catches errors at compile time, not runtime. When the backend expects `user_id: number` and the user accidentally sends a string, TypeScript immediately tells us where the error is. This saves hours of debugging, if we were to rely on JavaScript.

2.2. Backend: FastAPI

Why FastAPI and not Flask/Django? The answer is simple: Async. When 10 users upload documents simultaneously, Flask handles them one-by-one. FastAPI handles them concurrently. The difference is 10 seconds vs 2 seconds for batch operations.

2.3. Database: PostgreSQL + PGVector

Why one database instead of two? Simplicity. PGVector lets us store vectors right next to the document metadata. When we delete a document, the cascade removes both the metadata *and* its vectors automatically. This simplifies things a lot.

Why PostgreSQL specifically? PostgreSQL is distinguished by its maturity and unwavering reliability. Its resilience is most evident during system failures. In the event of a power loss or mid-process interruption, PostgreSQL utilizes a Write-Ahead Logging (WAL) mechanism to recover and maintain consistency

2.4. AI: Sentence Transformers + Google Gemini

Why Sentence Transformers? They're locally runnable, free, and produce embeddings good for our case. The `all-MiniLM-L6-v2` model gives us 384-dimensional vectors that capture semantic meaning without being massive.

Why Google Gemini? It's currently the best quality/cost ratio for the question-answering part.

3. Frontend Architecture

3.1. Folder Structure

The frontend follows a clear separation:

```
frontend/  
  src/  
    components/          # The reusable UI pieces part of code
```

```

Navbar.tsx
DocumentUpload.tsx
DocumentList.tsx
ChatInterface.tsx
pages/          # Full page components
  dashboard
    DashboardPage.tsx
  login
    LoginPage.tsx
  register
    RegisterPage.tsx
App.tsx          # Our root component
main.tsx         # The react bootstrap

```

Each component has one job. `Navbar` just shows the user and logout button. `DocumentUpload` handles file selection and upload. This makes debugging trivial.

3.2. State Management: Lifted State in App.tsx

In this case, we're not gonna use libraries like Redux nor Context API because it'd be overkill for this scale: we wouldn't need 'global state' data that would be needed by different components at once. Instead, I lift state up to `App.tsx`; I shared the states common to their parent:

```

1 // App.tsx - Central state management
2 export default function App() {
3   const [isLoggedIn, setIsLoggedIn] = useState(false);
4   const [username, setUsername] = useState('');
5   const [userId, setUserId] = useState(0);
6
7   // This state is passed down to all children
8   return (
9     <div className="app-container">
10       {isLoggedIn ? (
11         <DashboardPage
12           username={username}
13           userId={userId}
14           onLogout={handleLogout}
15         />
16       ) : (
17         <LoginPage onSuccess={handleAuthSuccess} />
18       )}
19     </div>
20   );
21 }

```

The `App.tsx` knows about authentication, then passes everything down as props. Each child component only gets what it needs.

3.3. UI Behavior During Async Operations

The UI never freezes; when the user uploads a document:

1. Button says "Uploading..." and disables
2. File input becomes disabled
3. The user can still navigate, chat, delete other documents
4. Progress could be added with axios interceptors

The key is `setUploading(true)` happening immediately, before the actual upload starts. This gives instant feedback.

3.4. File Upload Mechanism

The upload uses `FormData`, not JSON:

```
1 const formData = new FormData();
2 formData.append('file', file);
3 formData.append('user_id', userId.toString());
4
5 await axios.post('http://localhost:8000/documents/upload', formData, {
6   headers: { 'Content-Type': 'multipart/form-data' }
7 })
```

This is crucial because:

- Files can be large
- `FormData` streams the file. It doesn't load it all into memory
- The browser shows native progress for `FormData`
- It works with any file type without encoding issues

4. Backend Architecture

4.1. Entry Point: `main.py`

The `main.py` file is the front door. Every request goes here first:

```
1 app = FastAPI(
2     title="DocuMind AI API",
3     description="Made by Yassine",
4     version="1.0.0"
5 )
```

```

6 # Global error handling
7 @app.exception_handler(Exception)
8 async def global_exception_handler(request: Request, exc: Exception):
9     print(f"Uncaught exception: {exc}")
10    return JSONResponse(status_code=500, content={"detail": "Internal error"})
11
12 # CORS middleware
13 app.add_middleware(
14     CORSMiddleware,
15     allow_origins=["http://localhost:5173"],
16     allow_credentials=True,
17     allow_methods=["*"],
18     allow_headers=["*"],
19 )
20
21
22 # Include routers
23 app.include_router(auth.router, prefix="/auth", tags=["Authentication"])
24 app.include_router(documents.router, prefix="/documents", tags=[
25     "Documents"])

```

Notice the CORS configuration: it only allows the frontend origin. If there were a real world project, this would be a real domain.

4.2. Router Separation:

Routers are like departments in a company:

- auth.py handles login/register
- documents.py handles file operations
- Each has its own responsibility

It simplifies things: If authentication breaks, We look in auth.py. If file upload breaks, We look in documents.py.

4.3. Dependency Injection: get_db()

FastAPI's dependency injection:

```

1 def get_db():
2     db = SessionLocal()
3     try:
4         yield db
5     finally:
6         db.close()
7
8 @router.post("/upload")
9 async def upload_document(db: Session = Depends(get_db)):

```

```

10     # db is automatically provided
11     # automatically closed after function

```

Every route gets a fresh database session of its own. If the route crashes, the session closes automatically.

5. Authentication System

5.1. auth.py as the Receptionist

The `auth.py` file is more like a the building receptionist:

- Register first checks in with a "Can I get a keycard?"
- The login checks: "Here's my keycard, let me in". It verifies it and lets the user in.
- The logout is simply ending the session with a simple "I'm leaving, bye".

It's simple at this version. No email verification, no password reset, no OAuth.

Those are important, but (hopefully) they'd be in a v2 of this project,

5.2. Password Hashing with bcrypt

Passwords are *never* stored plaintext:

```

1 pwd_context = CryptContext(schemes=["bcrypt"], deprecated="auto")
2
3 def hash_password(password: str):
4     return pwd_context.hash(password)  # Adds salt automatically
5
6 def verify_password(plain_password: str, hashed_password: str):
7     return pwd_context.verify(plain_password, hashed_password)

```

Bcrypt simplifies password security by integrating salting and hashing into a single, automated process. It generates a unique, random salt for every user to prevent rainbow table attacks, ensuring identical passwords result in different hashes. The algorithm employs a "work factor" to perform multiple iterations, making brute-force attempts computationally expensive. To protect against timing attacks, it performs comparisons in constant time, preventing attackers from deducing secrets based on response speeds. Consequently, developers can implement high-level cryptographic security without managing low-level primitives manually.

5.3. Session Management

In this project, we used `withCredentials: true` in axios calls. This sends cookies automatically. The backend doesn't do JWT tokens because:

1. Cookies are simpler
2. They're automatically sent with every request
3. They have expiration built-in
4. The browser handles storage

For this version of this project, it is sufficient.

6. Database Design

6.1. The Database Defense

Data validation happens at three layers:

Layer	Technology	Example
Frontend	TypeScript	<code>userId: number type</code>
API	Pydantic	<code>constr(min_length=3)</code>
Database	SQL Constraints	<code>CHECK(LENGTH(username) > 0)</code>

If a malicious request somehow bypasses TypeScript and Pydantic, the database constraint catches it.

6.2. Entities Of Our Database:

Table	Purpose	Critical Fields
<code>users</code>	Store user accounts	<code>id, username, password_hash</code>
<code>documents</code>	Document metadata	<code>id, user_id, file_name, uploaded_at</code>
<code>document_chunks</code>	Text chunks + vectors	<code>id, document_id, user_id, content, embedding</code>

6.3. Foreign Keys with ON DELETE CASCADE

This is the most important design decision:

```

1 # In models.py
2 user_id = Column(Integer, ForeignKey("users.id", ondelete="CASCADE"))
3 document_id = Column(Integer, ForeignKey("documents.id", ondelete="CASCADE"))

```

When a user is deleted:

1. All their documents are deleted (cascade from users to documents)
2. All document chunks are deleted (cascade from documents to chunks)
3. No orphaned data
4. No manual cleanup code

The database handles all the integrity for us.

6.4. Data Isolation Per User

Every query includes `user_id`:

```

1 # Always filter by user_id
2 chunks = db.query(DocumentChunk).filter(
3     DocumentChunk.user_id == user_id
4 ).all()

```

User A never sees User B's documents, even if they guess the document ID.

7. AI & Document Processing Pipeline

7.1. From PDF Bytes to Text

When you upload a PDF, here's what happens:

```

1 def read_pdf(file_bytes):
2     doc = fitz.open(stream=file_bytes, filetype="pdf")
3     text = ""
4     for page in doc:
5         page_text = page.get_text()
6         text += page_text + "\n"
7     return text.strip()

```

`fitz` (the PyMuPDF library) reads the PDF directly from bytes, no temporary files. It handles:

- Encrypted PDFs, but with errors
- Scanned PDFs, but rather poorly, considering it's OCR territory)
- Multi-column layouts, rather adequately

7.2. Chunking Logic: The 500-Character Window

Why 500 characters? Too small (100 chars): we lose the context. Too large (1000 chars): embeddings get rather noisy. 500 is the char number we chose.

```

1 chunk_size = 500
2 chunks = []
3
4 for i in range(0, len(text), chunk_size):
5     chunk_text = text[i:i+chunk_size]
6     chunks.append(chunk_text)

```

Important: No overlap between chunks. Overlap helps only in doubling the storage.

7.3. Embedding Generation: Why 384 Dimensions

The `all-MiniLM-L6-v2` model produces 384-dimensional vectors. Each dimension represents some abstract feature of the text. Math example:

- "apple fruit" → [0.2, 0.8, -0.1, ..., 0.4]
- "orange citrus" → [0.3, 0.7, -0.2, ..., 0.3]
- "car vehicle" → [-0.8, 0.1, 0.9, ..., -0.5]

Vector embeddings transform words into multi-dimensional coordinates, allowing AI to interpret semantic relationships as mathematical distance. Words with shared contexts, such as "apple" and "orange," are mapped within similar spatial vectors, resulting in a low cosine distance (a measure that calculates the distance of high-dimensional vectors).

Conversely, unrelated terms like "apple" and "car" are positioned far apart in this vector space, resulting in a high distance. By calculating the angle between these vectors, the AI identifies topical relevance based purely on a distance value.

7.4. Singleton Pattern for Model Loading

Loading the embedding model takes: 2 seconds and 400MB RAM. We can't do that per request:

```

1 _embedding_model = None # Global variable
2
3 def get_embedding_model():
4     global _embedding_model
5     if _embedding_model is None:
6         _embedding_model = SentenceTransformer('all-MiniLM-L6-v2')
7     return _embedding_model

```

We utilize the Singleton Pattern: We ensure the model is instantiated only once and persisted in memory for all subsequent operations.

We then load the SentenceTransformer into a global variable; however, for the next thousand requests, we simply return its existing reference.

7.5. Semantic Search with Cosine Distance

Finding relevant chunks uses vector math:

```

1 # Convert question to vector
2 question_vector = model.encode(question)
3
4 # PostgreSQL finds closest vectors
5 results = db.query(DocumentChunk).filter(
6     DocumentChunk.user_id == user_id
7 ).order_by(
8     DocumentChunk.embedding.cosine_distance(question_vector)
9 ).limit(5).all()

```

The pgvector extension optimizes similarity searches by utilizing the Hierarchical Navigable Small World (HNSW) index, which organizes high-dimensional vectors into a multilayered graph. Rather than performing an exhaustive search across all records, the system traverses these layers to rapidly locate the nearest neighbors. This process relies on calculating the **Cosine Distance**, represented by the formula:

$$d = 1 - \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|}$$

By navigating this mathematical graph, the database can identify topically relevant content within milliseconds, even when scaling to millions of entries. This ensures that the retrieval process remains efficient without sacrificing the semantic accuracy of the results.

7.6. Question Answering with Gemini

The final step: send context + question to Gemini:

```

1 prompt = f"""
2 You are a helpful assistant. Answer based strictly on context.
3
4 CONTEXT:
5 {context_text}
6
7 QUESTION:
8 {question}
9 """
10
11 response = model.generate_content(prompt)
12 return response.text

```

We use the "strictly on context" part to prevent Gemini from using its general knowledge, which defeats the purpose of RAG in entirety.

8. API Routes in Detail

8.1. POST /documents/upload

1. Validate file type (PDF, TXT, Excel, CSV)
2. Read file content based on type
3. Create document record
4. Split text into 500-character chunks
5. Generate embedding for each chunk
6. Store chunks with vectors
7. Commit transaction

If step 6 fails, the entire transaction rolls back. No half-uploaded documents.

8.2. GET /documents/list/{user_id}

Returns metadata only, not content:

```

1  [
2    {
3      "id": 1,
4      "file_name": "financial_report.pdf",
5      "uploaded_at": "2025-01-15T10:30:00Z"
6    }
7 ]

```

Why not content? Privacy and bandwidth. The list might show 50 documents - sending all content would be megabytes.

8.3. DELETE /documents/{document_id}

```

1 @router.delete("/{document_id}")
2 def delete_one(document_id: int, db: Session = Depends(get_db)):
3     doc = db.query(Document).filter(Document.id == document_id).first()
4     if doc:
5         db.delete(doc) # Cascade deletes chunks automatically
6         db.commit()
7     return {"message": "Deleted"}

```

The ON DELETE CASCADE is pretty crucial: we delete the document, PostgreSQL deletes all its chunks.

8.4. DELETE /documents/all/{user_id}

This deletes every document for a user. A confirmation dialog in the frontend is included.

8.5. POST /documents/query

The RAG endpoint:

1. Generate question embedding
2. Find 5 most similar chunks
3. Concatenate as context
4. Send to Gemini
5. Return answer

If no documents exist, we simply return "Upload documents first" instead of crashing.

9. Infrastructure & Deployment

9.1. Docker Compose Orchestration

The docker-compose.yml defines three services:

```

1 services:
2   db:
3     image: ankane/pgvector:latest  # PostgreSQL + vector extension
4     environment:
5       POSTGRES_PASSWORD: 2040
6     volumes:
7       - postgres_data:/var/lib/postgresql/data
8
9   backend:
10    build: ./backend
11    ports: ["8000:8000"]
12    depends_on: [db]
13
14   frontend:
15    build: ./frontend
16    ports: ["5173:5173"]
17    depends_on: [backend]
```

One command: docker-compose up, and everything starts.

9.2. Startup Race Conditions

PostgreSQL takes 5-10 seconds to start. The backend would crash if it tried to connect immediately.

We therefore included this solution: The `startup.sh` file:

```

1 echo "Waiting for database to be ready."
2 sleep 5 # Simple but effective
3 python init_db.py
4 uvicorn app.main:app --host 0.0.0.0 --port 8000

```

9.3. Database Connection Retry Logic

In `database.py`:

```

1 def create_engine_with_retry(max_retries=10, delay=5):
2     for attempt in range(max_retries):
3         try:
4             engine = create_engine(DATABASE_URL)
5             with engine.connect() as conn:
6                 conn.execute(text("SELECT 1"))
7             return engine
8         except OperationalError:
9             if attempt == max_retries - 1:
10                 raise
11             time.sleep(delay)

```

If the database isn't ready, it waits 5 seconds and try again. This process is repeated 10 times. This handles slow startup and temporary network issues.

9.4. Frontend Docker Optimizations

The frontend Dockerfile has two stages:

1. Build stage: Install dependencies, compile TypeScript
2. Production stage: Copy only compiled files, use nginx

This keeps the production image small: no dev dependencies nor any source code.

10. End-to-End Flow: Upload to Query

10.1. Step 1: User Selects File

Frontend: User clicks "Choose File", selects '`report.pdf`' for example. React then updates the state.

10.2. Step 2: Frontend Uploads

```
1 const formData = new FormData();
2 formData.append('file', file);
3 formData.append('user_id', '123');
4 axios.post('http://localhost:8000/documents/upload', formData);
```

The browser streams the file.

10.3. Step 3: Backend Receives

FastAPI receives multipart form data. `UploadFile` gives us a file-like object.

10.4. Step 4: Read PDF

`fitz.open(stream=file_bytes)` gives us the advantage of reading directly from memory.

10.5. Step 5: Chunking

2000-character PDF becomes 4 chunks (500 chars each). Each chunk gets an embedding.

10.6. Step 6: Database Insert

Single transaction:

- Insert document record (ID: 456)
- Insert chunk 1 with embedding
- Insert chunk 2 with embedding
- Insert chunk 3 with embedding
- Insert chunk 4 with embedding
- Commit

PostgreSQL executes an automatic Rollback. This ensures that no document record, without its corresponding chunks, is ever persisted. The database simply reverts to its state prior to the start of the transaction

10.7. Step 7: User Asks Question

"What's the summary?" → Frontend sends to `/documents/query`.

10.8. Step 8: Semantic Search

1. Embed question: "What's the summary?" → vector Q
2. Database: find chunks with smallest cosine distance to Q
3. Returns top 5 chunks

10.9. Step 9: Answer Generation

Chunks concatenated as context. Then, it's sent to Gemini with the prompt. Afterwards, the response returns.

10.10. Step 10: Display

Frontend shows answer in chat bubble (the frontend). Conversation continues in the same manner.

11. Error Handling & Resilience

11.1. Frontend:

When backend is down:

```

1 try {
2   const response = await axios.get('http://localhost:8000/auth/me');
3   handleAuthSuccess(response.data);
4 } catch (error) {
5   // Not "NetworkError", just show login screen
6   setIsLoggedIn(false);
7 }
```

No error messages are actually shown. We just show the login form as if the user wasn't logged in.

11.2. Backend: Structured Errors

Every error returns consistent format:

```

1 {
2   "detail": "Could not read PDF",
3   "type": "ValueError"
4 }
```

11.3. Database: Connection Pooling

SQLAlchemy connection pool:

```
1 engine = create_engine(  
2     DATABASE_URL,  
3     pool_size=10,          # 10 persistent connections  
4     max_overflow=20,       # Can create 20 more if needed  
5     pool_pre_ping=True    # Check connection before use  
6 )
```

This handles 30 concurrent requests without creating 30 database connections.

12. Project Structure Recap

12.1. Root Infrastructure

The project follows a decoupled microservices architecture, managed by Docker Compose. The root directory contains orchestration files and global configuration:

```
.  
.gitignore  
docker-compose.yml  
Doc/  
    backend/  
    frontend/  
    public/
```

12.2. Backend Architecture (FastAPI)

The backend is organized into a modular FastAPI structure, separating the API logic from the data layer:

```
backend/  
    app/  
        routers/  
            auth.py  
            documents.py  
        database.py  
        main.py  
        models.py  
    init_db.py  
    requirements.txt  
    startup.sh
```

12.3. Frontend Architecture (React & TypeScript)

The frontend utilizes a feature-based organization. State is "lifted" to `App.tsx` to maintain a clean data flow between the dashboard and its child components:

```
frontend/
src/
  components/
    ChatInterface.tsx
    DocumentList.tsx
    DocumentUpload.tsx
    Navbar.tsx
  pages/
    dashboard/DashboardPage.tsx
    login/LoginPage.tsx
    register/RegisterPage.tsx
  App.tsx
  main.tsx
  App.css
vite.config.ts
Dockerfile
```

13. What Works Well

13.1. The RAG Pipeline

The core functionality - upload, chunk, embed, search, answer - works reliably. For a v1, it's surprisingly effective.

13.2. Docker Compose Setup

One command to run everything. Perfect for development, decent for production.

13.3. Type Safety

TypeScript + Pydantic catches so many errors. The confidence it provides is worth the extra typing.

13.4. Database Design

The three-table schema with cascading deletes is elegant. It handles all use cases without complexity.

14. Intentional Simplifications

14.1. No User Sessions Database

Sessions are stored in FastAPI's memory. This means:

- Restarting backend logs everyone out
- Multiple backend instances don't share sessions

14.2. No File Size Limits

A user could upload a 1GB PDF and crash the system.

14.3. No Rate Limiting

14.4. No Background Processing

Large PDFs block the upload endpoint.

15. Conclusion

DocuMind AI started as a simple question: “Can I make Ctrl+F smarter?” It became a full-stack application with React, FastAPI, PostgreSQL, and vector search.

The system works: You can upload documents, ask questions, get answers based on your documents.

Most importantly, I built every part myself. From the Docker containers to the SQL migrations to the React state management. I understand how it all fits together, where it breaks, and how to fix it.

The code is available on GitHub: <https://github.com/yassineexng/documind>. Feel free to use it, extend it, or learn from it.