

# Legal Document QnA System (RAG) - Project Documentation

---

## 1. Project Overview

**Project Name:** Legal Document QnA System (RAG)

**Duration:** 1 week

**Complexity Level:** Intermediate

**Description:** The Legal Document QnA System is a Retrieval-Augmented Generation (RAG) system designed for comprehensive legal document analysis. It allows users to upload large legal documents and query them via a chat interface. The system integrates advanced compression techniques to optimize embedding and query efficiency.

**Technical Stack:** - **Backend:** Python, FastAPI - **Vector Database:** ChromaDB (version 1.4.1) - **Frontend:** Next.js (default project structure, with upload + chat interface) - **Embeddings:** OpenAIEmbeddings - **LLM:** GPT-4o-mini / ChatOpenAI - **Optional Compression API:** ScaleDown

---

## 2. Key Features

Feature	Description	Benefit
Multi-Document Analysis	Users can upload multiple legal documents simultaneously.	Enables cross-document reasoning and comparison.
Citation Extraction	Extracts citations, case references, and legal sections from text.	Enhances answer credibility and traceability.
Precedent Linking	Links query answers to related precedents or documents.	Supports legal research and contextual understanding.
Confidence Scoring	Provides a relevance/confidence score for each answer.	Helps users assess reliability.
Compression with ScaleDown	Compresses large legal documents before embedding.	Reduces embedding costs by 60-80% and allows 5x larger docs in same context window.
Fast Query Response	Optimized retrieval and embeddings.	Reduces latency from 3-4s to under 1s.

---

## 3. Architecture Overview

### 3.1 Backend Pipeline

1. **Document Upload:** Accepts PDF, DOCX, TXT files.
2. **Preprocessing & Compression:** Uses ScaleDown API (optional) to compress text.
3. **Chunking:** Splits text into 800-token chunks with 150-token overlap.
4. **Embedding:** Generates vector embeddings via OpenAIEmbeddings.
5. **Vector Database Storage:** Stores embeddings in ChromaDB.
6. **QA Module:**
7. Retrieves relevant chunks based on query.
8. Combines context and generates answer with LLM.
9. Returns answer with source references and confidence.

### 3.2 Frontend

- Next.js chat interface.
- File upload component.
- Display for answer + sources + confidence score.
- Optional metrics dashboard showing token reduction from ScaleDown.

### 3.3 System Flow Diagram

```
User Upload -> Backend Preprocessing -> ScaleDown Compression -> Chunking ->
Embeddings -> ChromaDB -> Query -> Retrieve Relevant Chunks -> LLM Answer
Generation -> Return Answer + Sources
```

---

## 4. Technical Implementation

### 4.1 Backend

#### File Structure:

```
backend/
|
├── main.py          # FastAPI endpoints
├── ingest_pipeline.py # Document ingestion, compression, chunking, embedding
├── qa.py            # Question answering logic
├── chroma_utils.py  # ChromaDB helper functions
└── scaledown.py     # Optional ScaleDown integration
```

```
└── requirements.txt  
└── data/raw_docs/      # Sample documents
```

#### Key Functions:

- `ingest_pdf(pdf_path: str)`: Ingests a PDF, compresses, chunks, and embeds into ChromaDB.
- `compress_text_with_scaledown(text: str)`: Compresses text; safely falls back if API unreachable.
- `answer_question(question: str)`: Queries ChromaDB, retrieves relevant chunks, generates LLM answer.

## 4.2 Frontend (Next.js)

**Key Pages/Components:** - `pages/index.tsx`: Chat interface and upload button. - `components/ChatBox.tsx`: Handles user query input and displays responses. - `components/FileUploader.tsx`: Allows multi-file uploads.

#### API Integration:

```
const res = await fetch('http://localhost:8000/ask', {  
  method: 'POST',  
  headers: {'Content-Type': 'application/json'},  
  body: JSON.stringify({question: userQuestion})  
});  
const data = await res.json();
```

## 5. ScaleDown Compression

**Purpose:** - Reduce token length of documents before embedding. - Lower API costs. - Handle larger documents. - Reduce query latency.

**Metrics:** - Token reduction: 60-80% - Query latency: <1s - Document size handled: 5x larger than without compression

**Fallback Handling:** - If API unavailable, the system continues with uncompressed text. - Ensures robust ingestion pipeline.

## 6. Deliverables

1. **Working Web App:** Upload + chat interface with LLM answers.
2. **Backend API Endpoints:**

3. `POST /ask` : Accepts question, returns answer with sources.
  4. `GET /health` : Health check endpoint.
  5. **Vector DB:** ChromaDB with sample legal document embeddings.
  6. **Compression Metrics Dashboard:** Token reduction statistics per document.
  7. **Documentation:** This canvas file and optional README for setup instructions.
- 

## 7. Setup Instructions

### Backend

```
cd backend
python -m venv venv
venv\Scripts\activate
pip install -r requirements.txt
python ingest_pipeline.py # Run once to ingest sample docs
uvicorn main:app --reload
```

### Frontend

```
cd frontend
npm install
npm run dev
```

**Access:** `http://localhost:3000`

---

## 8. Unique Selling Points

- Handles **extremely large legal PDFs** with compression.
  - Multi-document retrieval + **precedent linking**.
  - Returns answers with **citations and confidence scores**.
  - Scalable vector database ingestion.
  - Robust fallback if external APIs are unavailable.
  - Fast query response (<1s) even on large datasets.
- 

## 9. Future Enhancements

- Fine-tune embeddings on legal corpus for better relevance.
- Multi-lingual support (English + Hindi).
- Interactive dashboard with query history and analytics.
- Support for offline querying with local LLMs.

---

**Project Complete:** Legal Document QnA System (RAG) with full ingestion, compression, retrieval, and interactive frontend.