# DocChat AI

**DocChat: Local Document Processing & AI Chat Platform**

DocChat is a local document processing platform that enables users to upload documents, extract text, and interact with a local AI chatbot powered by Ollama LLM. It preserves document formatting and provides a seamless workflow for both document handling and AI-assisted queries.
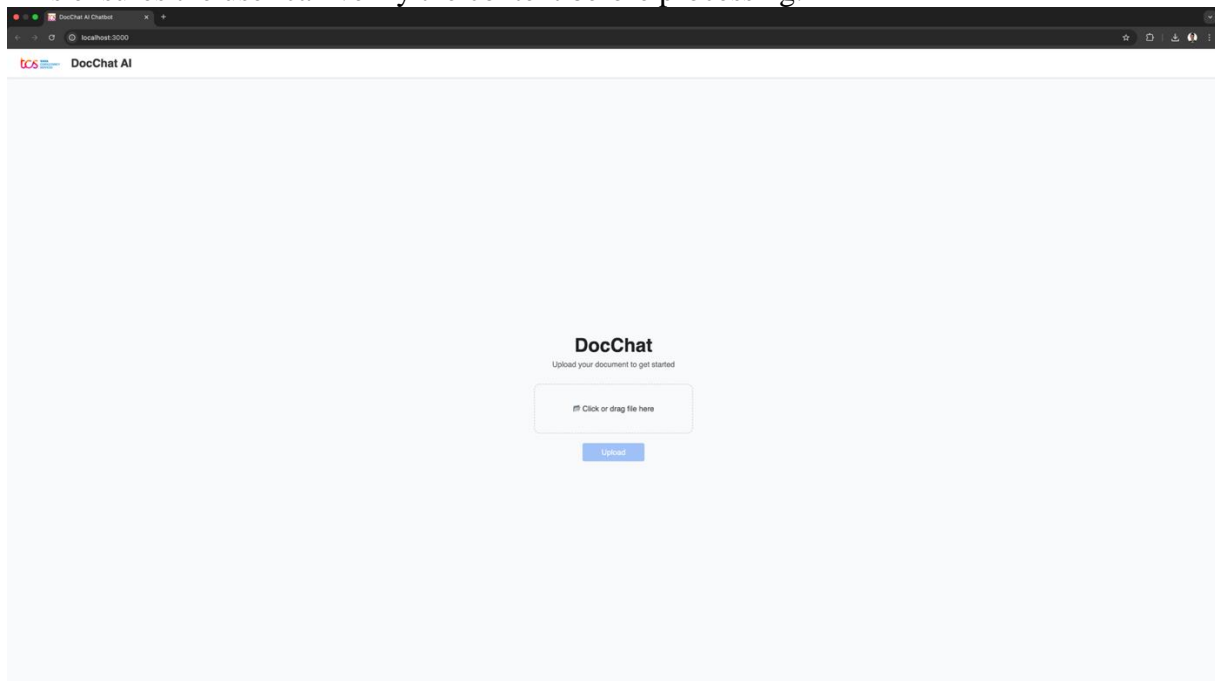
**Key Features:**

- Upload documents in .txt, .pdf, .docx, and .xlsx formats to local S3 storage.
- Extract text while maintaining original formatting.
- Chat with a local AI model (Ollama LLM) to summarize content or answer document-related questions.
- Fully local AI processing; no cloud AI calls are required.

**Tech Stack:**

- **Frontend:** React with file uploads, chat interface, and document previews. Axios handles API requests.
- **Backend:** Node.js + Express, providing REST APIs for document management and AI interactions. Integrates with local S3 (AWS SDK) and Ollama LLM.
- **Document Parsing:** pdf-parse-fixed, mammoth (DOCX), xlsx (Excel).
- **Local AI Model:** Ollama LLM (llama2) running at http://localhost:11434.

**Step 1: Upload Document**

- Users can upload any document: `.txt`, `.pdf`, `.docx`, `.xlsx`.
- After selecting a file, a preview is displayed on the right panel.
- This ensures the user can verify the content before processing.

| | GRADE: VIII A | UNIT TEST II MARKSHEET -2025-26 | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| S.NO | NAME OF THE STUDENT | MATH | SCIENCE | II LANG | ENG | SST | GRAND TOTAL | % | GRADE |
| | | 20 | 20 | 20 | 20 | 20 | 100 | | |
| 1 | Aarthi | 5 | 4 | 4 | 7 | 12 | 32 | 25.6 | E |
| 2 | Amirthavarshini | 0 | 8 | 1.5 | 8.5 | 11 | 29 | 23.2 | E |
| 3 | Deeksha | 10 | 19 | 16 | 17.5 | 17 | 79 | 63.2 | B2 |
| 4 | Gelling Laasya | ab | 19.5 | 15 | 18.5 | 13 | 66 | 52.8 | C1 |
| 5 | Manasa | 4 | 14.5 | 12 | 17.5 | 14.5 | 63 | 50 | C2 |
| 6 | Mani megalai | 5 | 1 | ab | 15 | | 21 | 16.4 | E |
| 7 | Navya Baghel | 7 | 18 | 19 | 17 | 13.5 | 75 | 59.6 | C1 |
| 8 | Priyadharshini | 17 | 18 | 19.5 | 18 | 20 | 93 | 74 | B1 |
| 9 | Ria Jeannie | 6 | ab | 8 | 19 | 9.5 | 43 | 34 | D |
| 10 | Rumaiza Zainab | 9 | 19.5 | 16 | 19.5 | 20 | 84 | 67.2 | B2 |
| 11 | Srinidhi | 10 | ab | 11.5 | 18 | 18 | 58 | 46 | C2 |
| 12 | Niharshia | ab | 7 | ab | 18 | 2.5 | 28 | 22 | E |
| 13 | Adhesh | 8 | 18.5 | 14 | 19 | 15.5 | 75 | 60 | C1 |
| 14 | Adithya | | 16.5 | 14 | 17 | 17 | 65 | 51.6 | C1 |
| 15 | Adhvaith Nair | 1 | 11.5 | 9 | 4.5 | | 26 | 20.8 | E |
| 16 | Denis Ray | 11 | 19 | 15 | 19 | 15 | 79 | 63.2 | B2 |
| 17 | Dilip | 3 | 5 | 2 | 2.5 | 8 | 21 | 16.4 | E |
| 18 | Fadil | 16 | 19 | 7.5 | 20 | 19 | 82 | 65.2 | B2 |
| 19 | Gautham | 3 | 15 | 10 | 17 | 18.5 | 64 | 50.8 | C1 |
| 20 | Harshith | 2 | 6.5 | 6.5 | 19.5 | 10.5 | 45 | 36 | D |
| 21 | Kavin Karthik | 5 | 6 | 3 | 5 | 3 | 22 | 17.6 | E |
| 22 | Mithilaesh | 7 | 9.5 | 5 | 12.5 | 16 | 50 | 40 | D |
| 23 | Monesh Kumar | 11 | 20 | 13.5 | 19.5 | 18.5 | 83 | 66 | B2 |
| 24 | Nalankili Chozhan | 11 | 17.5 | 18 | 19.5 | ab | 66 | 52.8 | C1 |
| 25 | Prajith | 8 | 19.5 | 16 | 18 | 16 | 77.5 | 62 | B2 |
| 26 | Pranith | 2 | ab | 7.5 | 16.5 | 9 | 35 | 28 | E |
| 27 | Richard | 2 | ab | 10 | 15.5 | 9 | 36.5 | 29.2 | E |
| 28 | Rishi | 2 | 16 | | 11 | 11.5 | 52.5 | 42 | C2 |
| 29 | Shaik Ayaz | 3 | 17.5 | 13.5 | 13.5 | 13.5 | 61 | 48.8 | C2 |

## Step 2: Upload to S3

- Once uploaded, the document is stored in the configured S3 bucket (or localstack for testing).
- The uploaded document can be verified in the S3 storage.
- This ensures secure storage and easy retrieval.

**Step 3: View & Chat Page**

- The document is displayed on the left panel for easy reference.
- The right panel hosts the AI chat interface.



**Step 4: Summarize Document**

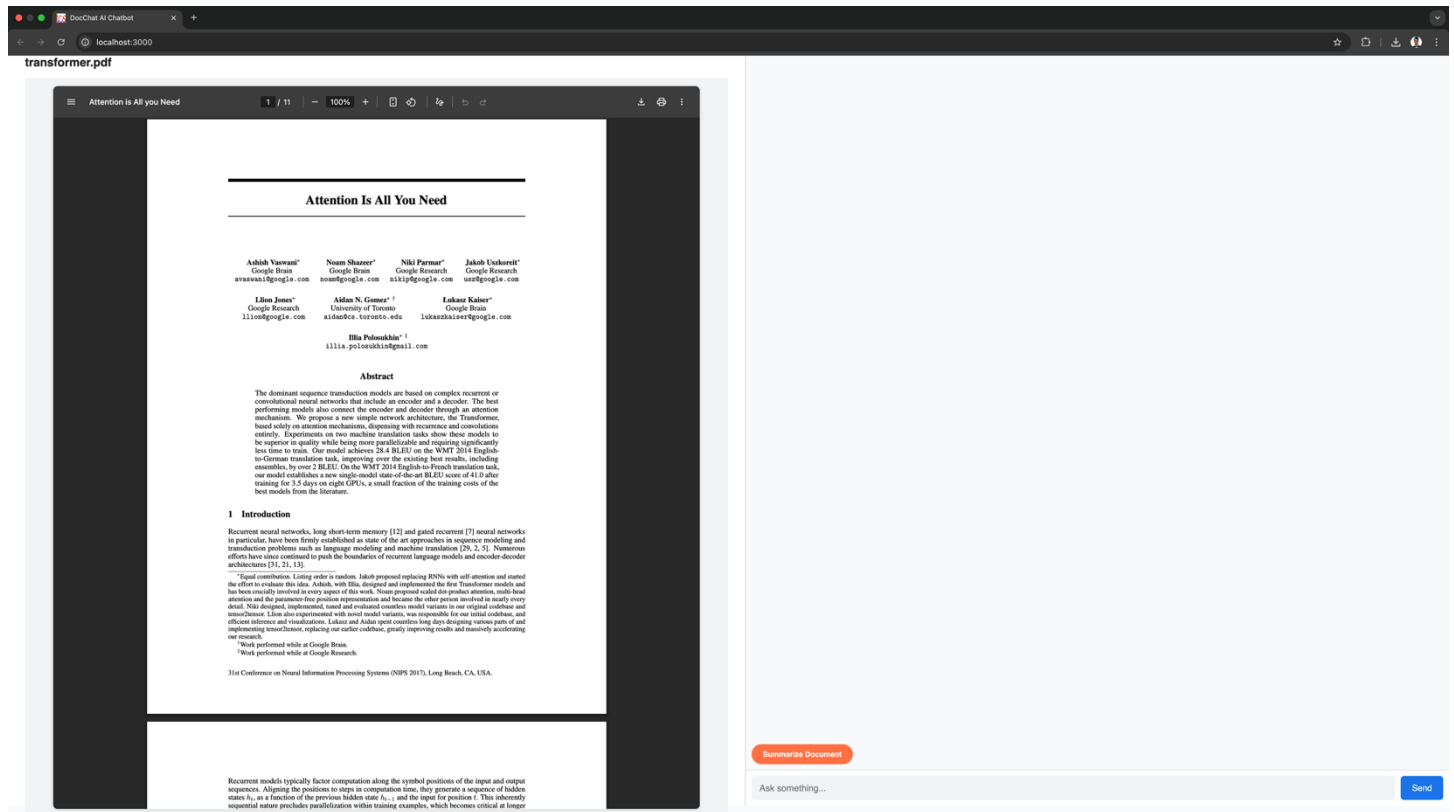- The **Summarize** button generates a concise summary of the document.
- This provides a quick overview without reading the entire content.

**Step 5: Chat with Local AI**

- Users can interact with the AI to ask questions or get detailed explanations.
- The AI is powered by **Ollama LLM**, running entirely **locally**, ensuring **full data privacy**.
- No cloud AI calls are required, keeping your documents secure.

# Benefits to the Enterprise

- **Increased Productivity:** Instant summaries and direct answers reduce manual scanning.
- **Faster Decision-Making:** Quick insights from complex reports or spreadsheets.
- **Enhanced Compliance:** Local AI ensures adherence to privacy and regulatory mandates.
- **Scalability:** Modular design handles thousands of documents and queries.
- **Competitive Advantage:** Differentiates the organization with secure, AI-driven document intelligence.

# Future Scope

- Additional document types (PowerPoint, emails, OCR for scanned documents).
- Multilingual document summarization and querying.
- Domain-specific LLM fine-tuning for finance, healthcare, or legal content.
- Voice interface integration.
- Advanced analytics for tabular data and trend forecasting.
- Knowledge graph integration to link insights across multiple documents.

# Conclusion

- DocChat AI represents a next-generation enterprise document assistant, integrating intelligence directly into document workflows. By leveraging locally hosted AI, secure parsing pipelines, and interactive conversational interfaces, it provides enhanced productivity, faster decision-making, and full compliance with enterprise security requirements.
- The solution is scalable, reusable, and adaptable, enabling **clients such as TCS** to adopt secure AI-driven document services across multiple industries and business domains.