

Building a Multi-Document RAG System



AMAN KHARWAL · JANUARY 6, 2026

AmanXai



By Aman Kharwal

Build a Multi-Document RAG System.



RAG connects the powerful reasoning of an LLM with the unique information in your own documents. Today, I'll teach you how to build a Multi-Document RAG System using Python. By the end, you'll have an app that reads a folder of documents and answers your questions accurately.



■ Multi-Document RAG System: Getting Started

We are going to build a Multi-Document RAG system from scratch using Python, LangChain, and Ollama. It sounds complex, but I promise you, it's just a series of logical steps.

We'll use LangChain for orchestration, Chroma for storage, and Ollama to run the Llama 3 model locally.

First, install these libraries. In your terminal, run:

```
pip install langchain langchain-community langchain-huggingface langchain-chroma langchain-ollama pypdf
```

You'll also need [Ollama](#) running locally with the Llama 3 model. After installing Ollama, run **ollama pull llama3**.

Step 1: Loading the Raw Knowledge



First, gather your source materials. We need to extract text from PDF files. PyPDFLoader is a good choice because it handles the tricky formatting of PDFs well:

```
1 import os
2 from langchain_community.document_loaders import PyPDFLoader
3
4 def load_documents(folder_path: str):
5     if not os.path.exists(folder_path):
6         raise FileNotFoundError(f"Folder '{folder_path}' does not exist")
7
8     documents = []
9     for filename in os.listdir(folder_path):
10         if filename.endswith(".pdf"):
11             file_path = os.path.join(folder_path, filename)
12             print(f"📄 Loading: {filename}")
13             try:
14                 loader = PyPDFLoader(file_path)
15                 documents.extend(loader.load())
16             except Exception as e:
17                 print(f"❌ Error loading {filename}: {e}")
18     return documents
```

- Data is rarely perfect. Make sure your loading logic skips non-PDFs and handles errors, so your pipeline keeps running even if one file is bad.

Step 2: Chunking

You can't give a 100-page document to an LLM all at once because it goes over the memory limit. So, we need to break it into smaller parts.

We use `RecursiveCharacterTextSplitter`, which is a smart tool. It tries to split text by paragraphs first, then by sentences, so related text stays together:

```
1 from langchain_text_splitters import RecursiveCharacterTextSplitter
2
3 def split_text(documents):
4     splitter = RecursiveCharacterTextSplitter(
5         chunk_size=1000,
6         chunk_overlap=200,
7     )
8     chunks = splitter.split_documents(documents)
9     print(f"✂ Created {len(chunks)} chunks")
```

```
10     return chunks
```

Pay attention to **chunk_overlap=200**. This setting is important because it creates a sliding window, making sure you don't lose context if a sentence is split between chunks.

Step 3: Embeddings

Computers understand numbers, not words. So, we need to turn our text chunks into lists of numbers, called **vectors** or **embeddings**.

This means that if two chunks have similar meanings, like "Dog" and "Puppy," their numbers will be close to each other:

```
1 from langchain_huggingface import HuggingFaceEmbeddings
2
3 embedding_function = HuggingFaceEmbeddings(
4     model_name="sentence-transformers/all-MiniLM-L6-v2"
5 )
```



■ We'll use all-MiniLM-L6-v2, a lightweight, open-source model that runs quickly on your CPU.

Step 4: The Vector Store

Now we need somewhere to store these numbers for fast searching. A regular SQL database isn't good for this, so we'll use a **Vector Database**. We'll use **Chroma**:

```
1 from langchain_chroma import Chroma
2
3 def create_vector_store(chunks):
4     vector_store = Chroma.from_documents(
5         documents=chunks,
6         embedding=embedding_function,
7         persist_directory="./chroma_db",
8         collection_name="rag_docs"
9     )
10     return vector_store
```

This function saves the database in a folder called **./chroma_db**. That way, you don't have to rebuild the database every time you restart the app; it stays saved.



■ Step 5: The Brain

This is the most important part. This function links the user, the database, and the LLM:

```
1 from langchain_ollama import ChatOllama
2 from langchain_core.prompts import ChatPromptTemplate
3 from langchain_core.runnables import RunnablePassthrough
4 from langchain_core.output_parsers import StrOutputParser
5
6
7 def format_docs(docs):
8     return "\n\n".join(doc.page_content for doc in docs)
9
10
11 def query_rag_system(query_text, vector_store):
12     llm = ChatOllama(model="llama3") # Make sure you have Ollama insta
13
14     retriever = vector_store.as_retriever(search_kwargs={"k": 3})
15
16     prompt = ChatPromptTemplate.from_template(
17         """
18         You are a helpful assistant.
```




```
19     Answer ONLY using the context below.
20     If the answer is not present, say "I don't know."
21
22     Context:
23     {context}
24
25     Question:
26     {question}
27     """
28 )
29
30 chain = (
31     {
32         "context": retriever | format_docs,
33         "question": RunnablePassthrough(),
34     }
35     | prompt
36     | llm
37     | StrOutputParser()
38 )
39
40 return chain.invoke(query_text)
```

First, it looks at the user's question and finds the top 3 most relevant chunks (**k=3**). Then, it puts those chunks into a strict prompt: **"Answer ONLY using the context below."** This helps stop the AI from making things up.

Step 6: Putting It All Together

Finally, the main loop checks if a database already exists. If it doesn't, it processes the PDFs. Then, it starts a chat loop so you can ask questions:

```
1 def main():
2     folder_path = "/Users/amankharwal/aiagent/data" # CHANGE THIS to y
3
4     if not os.path.exists("./chroma_db"):
5         print("📦 No vector DB found. Creating one...")
6         docs = load_documents(folder_path)
7         chunks = split_text(docs)
8         vector_store = create_vector_store(chunks)
9         print("Vector database created")
10    else:
11        print("📦 Loading existing vector DB...")
```



```
12     vector_store = Chroma(  
13         persist_directory="./chroma_db",  
14         embedding_function=embedding_function,  
15         collection_name="rag_docs"  
16     )  
17  
18     while True:  
19         query = input("\n ? Ask a question (or type 'exit'): ")  
20         if query.lower() == "exit":  
21             break  
22  
23         print("😬 Thinking...")  
24         answer = query_rag_system(query, vector_store)  
25         print("\n🧠 Answer:\n", answer)  
26  
27 if __name__ == "__main__":  
28     main()
```

Here's the answer I got for my and my friends' resumes:

```
(env) (base) amankharwal@Amans-MacBook-Pro aiagent % python multirag.py
```

```
📦 Loading existing vector DB...
```

```
? Ask a question (or type 'exit'): Compare Tushar and Aman for the role of a Data Analyst
```

```
🤔 Thinking...
```

```
🧠 Answer:
```

```
Based on the provided context, here's a comparison between Tushar and Aman for the role of a Data Analyst:
```

```
**Similarities:**
```

- * Both Tushar and Aman have experience with data-related tools like SQL, Power BI, and Excel.
- * They both have a strong foundation in analytical thinking.

```
**Differences:**
```

- * **Experience:** Aman has 5+ years of experience building intelligent systems using Machine Learning, NLP, and Generative AI, whereas Tushar is a recent graduate with limited work experience (February 2024 – September 2024).
- * **Specialization:** Aman specializes in developing LLM-based systems, RAG chatbots, and Multi-Agent architectures that solve real-world business challenges. Tushar has not mentioned any specific areas of specialization.
- * **Background:** Aman is the founder of Statso.io and has led data-driven projects for startups and enterprises, whereas Tushar is a recent graduate with limited work experience.

```
**Conclusion:**
```

Based on the provided context, Aman appears to be a more suitable candidate for the role of a Data Analyst due to his extensive experience in building intelligent systems using Machine Learning, NLP, and Generative AI. However, it's essential to consider Tushar's skills and background as well. If you're looking for a recent graduate with strong foundations in SQL, Power BI, and analytical thinking, Tushar might be a good fit.

Please note that this comparison is based solely on the provided context, and additional information may be necessary to make a more informed decision.

Closing Thoughts

Building systems like this shows me that AI isn't meant to replace our curiosity; it helps fuel it. When it's easier to find answers, we can ask better, deeper, and more creative questions.



- Don't be afraid to experiment with this code. Try changing the chunk size, swap llama3 for Mistral, or use a different embedding model. That's the best way to learn.

If you found this article useful, you can follow me on [Instagram](#) for daily AI tips and practical resources. You might also like my latest book, [Hands-On GenAI, LLMs & AI Agents](#). It's a step-by-step guide to help you get ready for jobs in today's AI field.



Aman Kharwal

AI/ML Engineer | Published Author. My aim is to decode data science for the real world in the most simple words.



ARTICLES: 2046





PREVIOUS POST

The Agentic AI Engineer Roadmap for 2026

NEXT POST

Top Companies Hiring for GenAI Roles

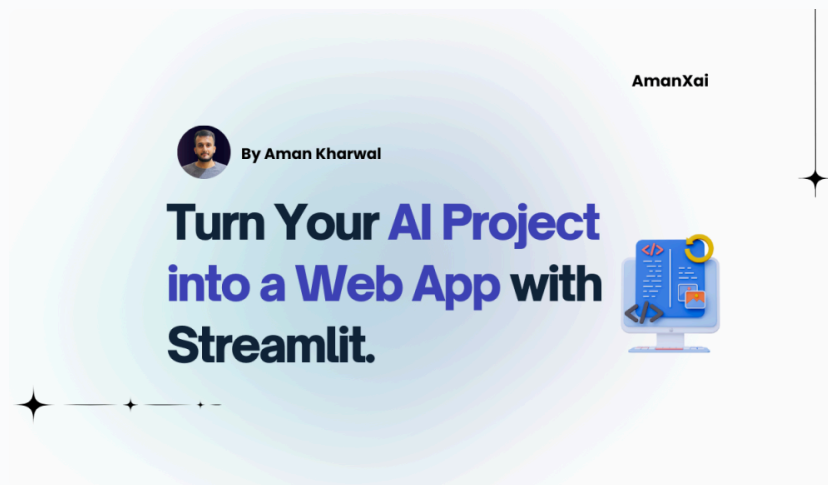


Recommended For You



4 End-to-End AI/ML Projects to Get Job-Ready

January 10, 2026



Build Your First AI App UI with Streamlit

January 8, 2026 / 3 Comments





AmanXai



By Aman Kharwal

A Complete List of Companies Hiring for GenAI Roles.



Top Companies Hiring for GenAI Roles

January 7, 2026

AmanXai



By Aman Kharwal

The **Agentic AI** Engineer Roadmap for 2026.



The Agentic AI Engineer Roadmap for 2026

January 5, 2026 / 2 Comments

Leave a Reply



