

PROJECTS

• The Daily News

• The Daily News

• The Daily News

• The Daily News

Retrieval Augmented Gen

Retrieval-Augmented Generation (RAG) is an artificial intelligence framework that enhances the performance of large language models (LLMs) by integrating them with external knowledge sources. Instead of relying solely on data the model was trained on, RAG retrieves relevant and up-to-date information from databases, documents, or other knowledge bases and incorporates this information into its generated responses

This approach addresses common limitations of LLMs, such as outdated knowledge and the risk of generating inaccurate or generic answers (often called "hallucinations"). By grounding responses in authoritative, current data, RAG improves accuracy, relevance, and transparency.

IITB, Mumbai

```
1 def rag_chatbot(query, db):
2     docs = retrieve_docs(query, db)
3     context = "\n\n".join([doc.page_content for doc in docs])
4     answer = generate_answer_groq(query, context)
5     return answer
6
```

RAG Q&A SYSTEM

```
1 import streamlit as st
2 import requests
3
4 import requests
5
6 GROQ_API_KEY = "i removed mine for security reasons, please set your own"
7
8 def generate_answer_groq(query, context, model="llama3-8b-8192"):
9     url = "https://api.groq.com/openai/v1/chat/completions"
10    headers = {
11        "Authorization": f"Bearer {GROQ_API_KEY}",
12        "Content-Type": "application/json"
13    }
14
15    system_prompt = "You are a helpful assistant that answers questions based only on the provided context."
16    user_prompt = f"Context:\n{context}\n\nQuestion: {query}\n\nAnswer:"
17
18    payload = {
19        "model": model,
20        "messages": [
21            {"role": "system", "content": system_prompt},
22            {"role": "user", "content": user_prompt}
23        ],
24        "temperature": 0.2,
25        "max_tokens": 300
26    }
27
28    response = requests.post(url, headers=headers, json=payload)
29
30    try:
31        data = response.json()
32        if "choices" in data:
33            return data["choices"][0]["message"]["content"]
34        else:
35            return f"⚠ Unexpected response: {data}"
36    except Exception as e:
37        return f"❌ Error: {e} | Raw response: {response.text}"
38
39
```

Classrooms of Sunset Valley School,

This book demonstrates a Retrieval-Augmented Generation (RAG) system designed to extract and answer questions from PDF documents. The workflow starts by loading and splitting the PDF text into overlapping chunks using PyPDF2 and LangChain's RecursiveCharacterTextSplitter, which helps retain context. These text segments are then further

indexed with FAISS and MiniLM embeddings, enabling efficient similarity searches. When a user submits a query, the system retrieves the most relevant document sections and generates accurate, context-aware answers using Groq's Llama3 API, which is limited to using only the provided context. An example shows the system successfully answering technical questions from an assignment PDF.

