

The learnings regarding RAG(Retrieval-Augmented Generation)

RAG (Retrieval-Augmented Generation) is a technique in AI and NLP (Natural Language Processing) that combines **information retrieval** and **text generation** to improve the quality and relevance of the generated output.

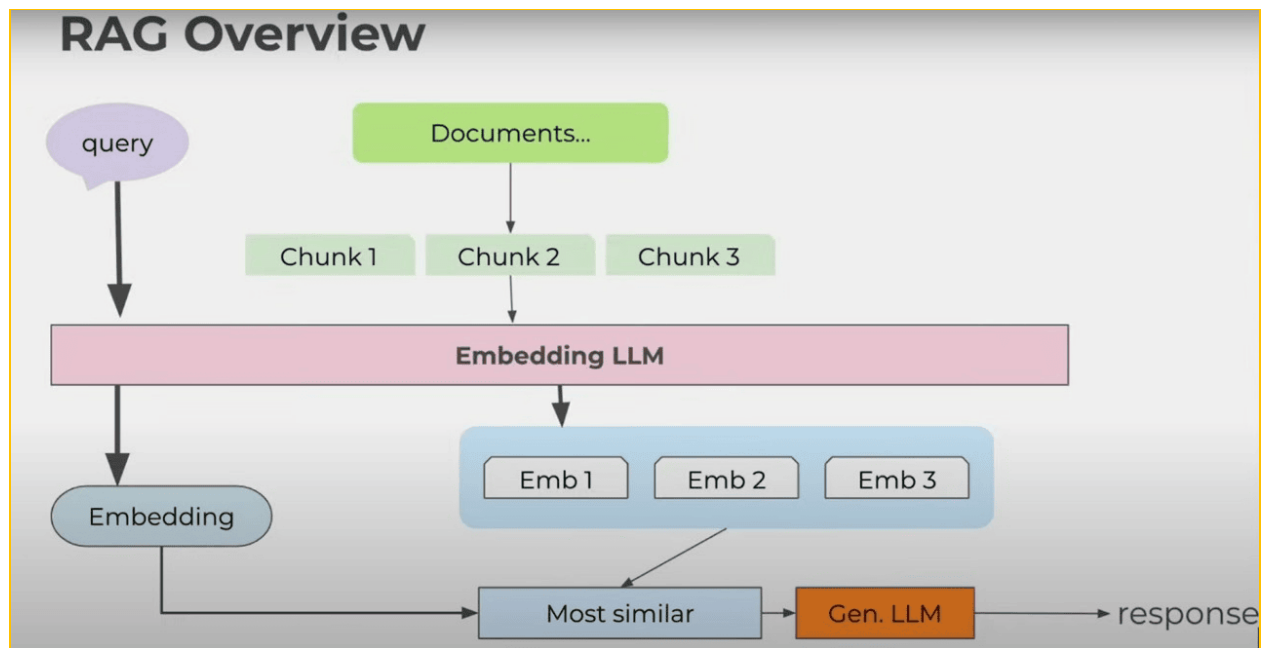


fig1: RAG overview

Steps in RAG (Retrieval-Augmented Generation):

- Data collection
- Data chunking
- Document embeddings
- Handling user queries
- Generating responses with an LLMs

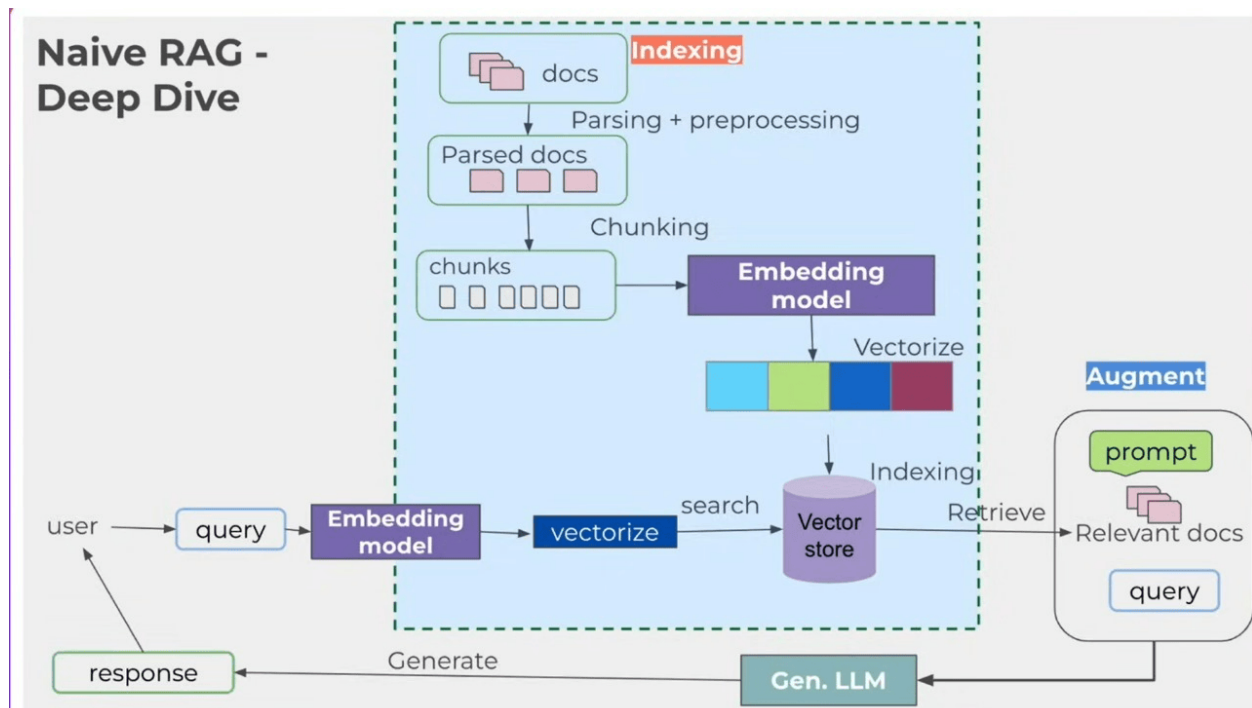


fig2: visualizing the steps

As a group project to complete the entire model of RAG, me and Vishal Pk were given the task of completing the first 2 steps i.e. data collection and data chunking.

We thought of making the RAG model as a medical chatbot which generates a response using gpt 2.0.

Data collection was done by using ChatGPT as we were not able to get the correct type of text for our project. we wanted the dataset to have the name of a specific disease and the discription in a ordered manner.

As we couldn't get that in huggingface , we tried adding the data manually by collecting from different sources which I have attached in references.

Data chunking was the next step. I have created a separate file that explains all the chunking methods and what it actually means in this workplace called [chunking docs\(1\)](#).

the codes which were considered what not used are:

```
1 from datasets import load_dataset
2
3 ds =
4     load_dataset("Somraj77/formatted_ai_medical_chatbot_llama2"
5     )
6
7 print("Available fields:", ds['train'].column_names)
8
9 prompts = []
10 responses = []
11
12 if 'text' in ds['train'].column_names:
13     texts = ds['train']['text']
14
15     for i, text in enumerate(texts, 1):
16
17         if "[INST]" in text and "[/INST]" in text:
18
19             prompt =
20                 text.split("[/INST]")[0].replace("[INST]", "").strip()
21             response = text.split("[/INST]")[1].strip()
22
23             prompt = prompt.replace("<s>",
24                 "").replace("</s>", "").strip()
25             response = response.replace("<s>",
26                 "").replace("</s>", "").strip()
27
28             prompts.append(prompt)
29             responses.append(response)
```

```

25         else:
26
27             print(f"Record {i} does not have proper markers
                and was skipped.")
28 else:
29     print("The dataset does not have a 'text' field. Check
        the available fields and use the correct name.")
30
31 print("\nPrompts List:")
32 print(prompts[:5])
33 print("\nResponses List:")
34 print(responses[:5])

```

this was first considered for a specific type of document, but we ended up using another one.

```

1  from PyPDF2 import PdfReader
2  import re
3
4
5  pdf_path = 'nihms-1718244.pdf'
6
7  reader = PdfReader(pdf_path)
8
9  extracted_text = ""
10 for page in reader.pages:
11     extracted_text += page.extract_text()
12
13
14 sentences = re.split(r'(?<=[.!?])\s+', extracted_text)
15
16

```

```

17 chunk_size = 1000
18 chunks = []
19 current_chunk = ""
20
21 for sentence in sentences:
22     if len(current_chunk) + len(sentence) + 1 <=
        chunk_size:
23         current_chunk += sentence + " "
24     else:
25         chunks.append(current_chunk.strip())
26         current_chunk = sentence + " "
27
28 if current_chunk:
29     chunks.append(current_chunk.strip())
30
31 output_path = 'chunked_text_for_RAG.txt'
32 with open(output_path, 'w', encoding='utf-8') as file:
33     for idx, chunk in enumerate(chunks, 1):
34         file.write(f"=== Chunk {idx} ===\n{chunk}\n\n")
35
36 print(f"Chunks saved to {output_path}")
37

```

this was not used for the same reasons.

```

1 chunks = []
2 with open('disease_dataset.txt', 'r') as file:
3     current_chunk = ""
4     for line in file:
5         if line.strip().startswith("Disease:") and
            current_chunk:
6

```

```

7         chunks.append(current_chunk.strip())
8         current_chunk = line
9     else:
10
11         current_chunk += line
12
13     if current_chunk:
14         chunks.append(current_chunk.strip())
15
16
17 for i, chunk in enumerate(chunks):
18     print(f"Chunk {i+1}:\n{chunk}\n")
19

```

the code which was used was:

```

1 def chunk_text_file(file_path, chunk_size):
2
3     chunks = []
4     with open(file_path, 'r', encoding='utf-8') as file:
5         current_chunk = []
6         for line in file:
7             if line.strip():
8                 current_chunk.append(line.strip())
9                 if len(current_chunk) == chunk_size:
10                     chunks.append("\n".join(current_chunk))
11                     current_chunk = []
12
13         if current_chunk:
14             chunks.append("\n".join(current_chunk))
15     return chunks
16
17
18 file_path = "input.txt"
19 chunk_size = 1

```

```
20 chunks = chunk_text_file("dataset.txt", chunk_size)
21
22
23
24 print(chunks)
25
```

This was used in the main code.

REFERENCES:

<https://www.datacamp.com/blog/what-is-retrieval-augmented-generation-rag>

<https://writer.zoho.in/writer/open/g0keb0e096a2d800a4201b3c9732aa0d4e785>

<https://docs.cohere.com/v2/docs/retrieval-augmented-generation-rag>

<https://rahuld3eora.medium.com/from-basic-to-advanced-rag-every-step-of-the-way-dee3a3a1aae9>

<https://docs.databricks.com/en/generative-ai/tutorials/ai-cookbook/fundamentals-data-pipeline-steps.html>

