

Retrieval Augmented Generation



A Simple Introduction

Abhinav Kimothi

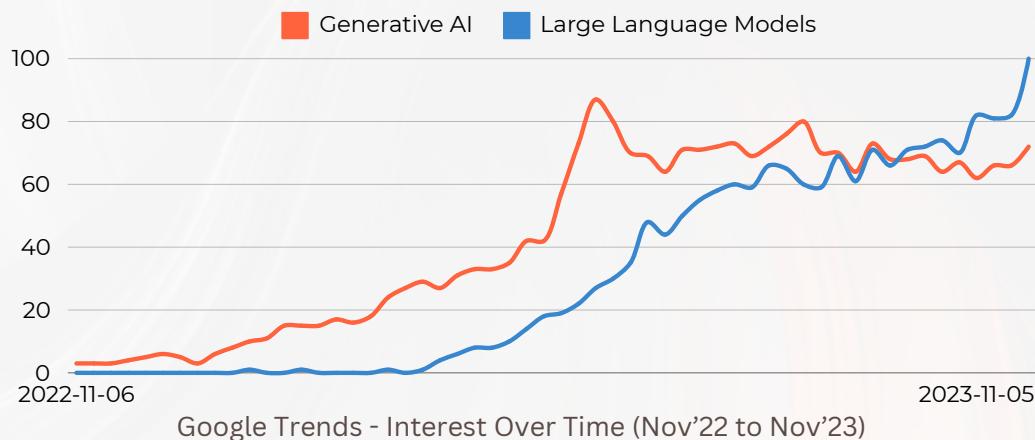
Table of Contents

01. What is RAG?	<u>3</u>
02. How does RAG help?	<u>6</u>
03. What are some popular RAG use cases?	<u>7</u>
04. RAG Architecture	<u>8</u>
i) Indexing Pipeline	<u>9</u>
a) Data Loading	<u>10</u>
b) Document Splitting	<u>14</u>
c) Embedding	<u>23</u>
d) Vector Stores	<u>29</u>
ii) RAG Pipeline	<u>35</u>
a) Retrieval	<u>37</u>
b) Augmentation and Generation	<u>45</u>
05. Evaluation	<u>46</u>
06. RAG vs Finetuning	<u>56</u>
07. Evolving RAG LLMOps Stack	<u>59</u>
08. Multimodal RAG	<u>63</u>
09. Progression of RAG Systems	<u>66</u>
i) Naive RAG	<u>66</u>
ii) Advanced RAG	<u>67</u>
iii) Multimodal RAG	<u>71</u>
10. Acknowledgements	<u>73</u>
11. Resources	<u>74</u>

What is RAG?

Retrieval Augmented Generation

30th November, 2022 will be remembered as the watershed moment in artificial intelligence. OpenAI released ChatGPT and the world was mesmerised. Interest in previously obscure terms like **Generative AI** and **Large Language Models (LLMs)**, was unstoppable over the following 12 months.



The Curse Of The LLMs

As usage exploded, so did the expectations. Many users started using ChatGPT as a source of information, like an **alternative to Google**. As a result, they also started encountering prominent weaknesses of the system. Concerns around copyright, privacy, security, ability to do mathematical calculations etc. aside, people realised that there are **two major limitations** of Large Language Models.

A Knowledge Cut-off date

Training an LLM is an expensive and time-consuming process. LLMs are trained on massive amount of data. The data that LLMs are trained on is therefore historical (or dated).

e.g. The latest GPT4 model by OpenAI has knowledge only till April 2023 and any event that happened post that date, the information is not available to the model.

Hallucinations

Often, it was observed that LLMs provided responses that were factually incorrect. Despite being factually incorrect, the LLM responses “sounded” extremely confident and legitimate. This characteristic of “lying with confidence” proved to be one of the biggest criticisms of ChatGPT and LLM techniques, in general.

Users look at LLMs for knowledge and wisdom, yet LLMs are sophisticated predictors of what word comes next.

The Hunger For More

While the weaknesses of LLMs were being discussed, a parallel discourse around providing context to the models started. In essence, it meant creating a ChatGPT on proprietary data.

The Challenge

- Make LLMs respond with up-to-date information
 - Make LLMs not respond with factually inaccurate information
 - Make LLMs aware of proprietary information
- Providing LLMs with information not in their memory

Providing Context

While model re-training/fine-tuning/reinforcement learning are options that solve the aforementioned challenges, these approaches are time-consuming and costly. In majority of the use-case, these costs are prohibitive.

In May 2020, researchers in their paper [“Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks”](#) explored models which combine pre-trained parametric and non-parametric memory for language generation.

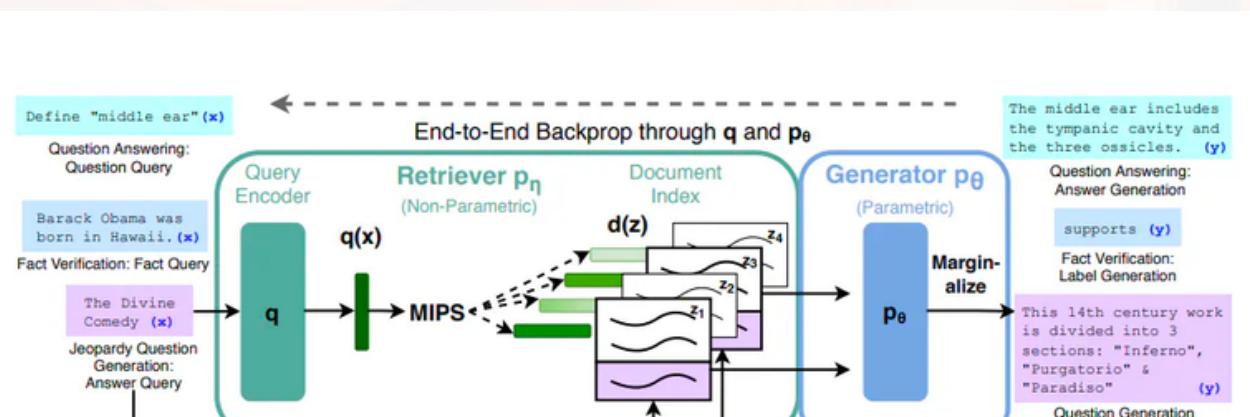
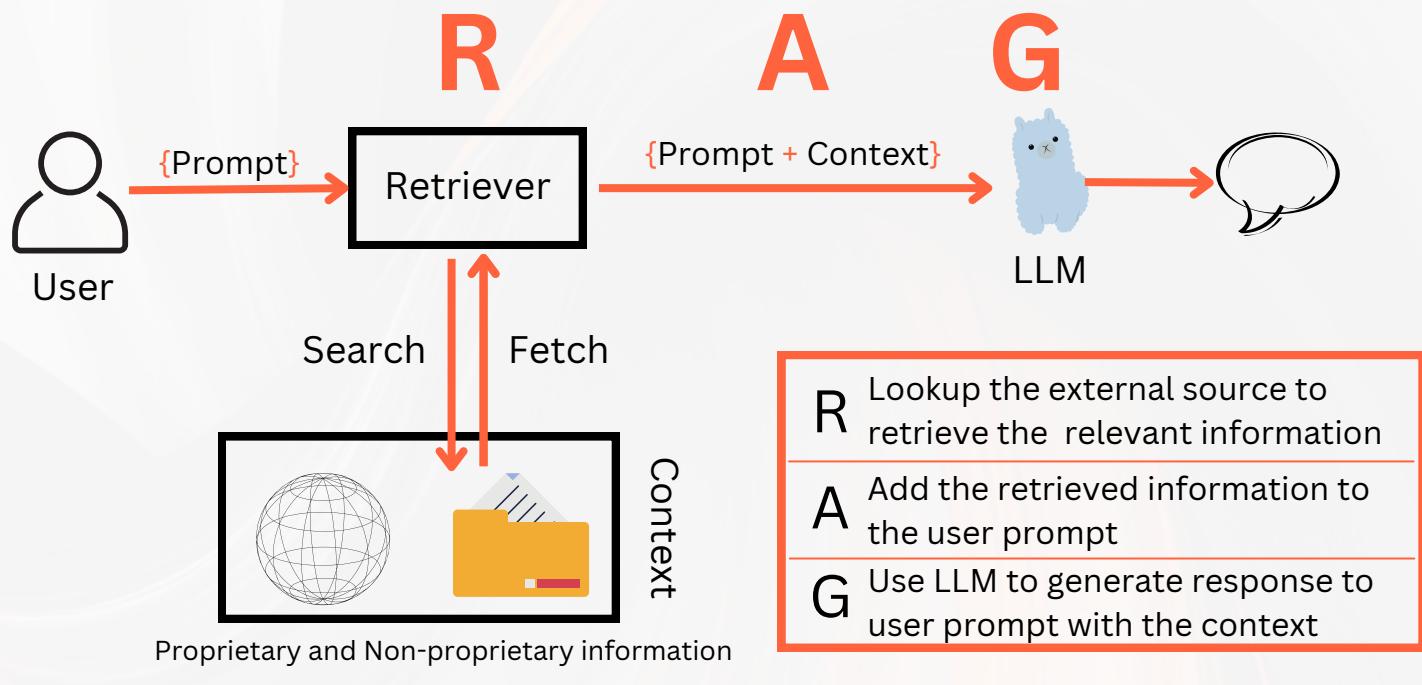


Figure 1: Overview of our approach. We combine a pre-trained retriever (*Query Encoder + Document Index*) with a pre-trained seq2seq model (*Generator*) and fine-tune end-to-end. For query x , we use Maximum Inner Product Search (MIPS) to find the top-K documents z_i . For final prediction y , we treat z as a latent variable and marginalize over seq2seq predictions given different documents.

So, What is RAG?

In 2023, RAG has become one of the most used technique in the domain of Large Language Models.



User enters a prompt/query

Retriever searches and fetches information relevant to the prompt (e.g. from the internet or internet data warehouse)

Retrieved relevant information is augmented to the prompt as context

LLM is asked to generate response to the prompt in the context (augmented information)

User receives the response

A Naive RAG workflow

How does RAG help?

Unlimited Knowledge

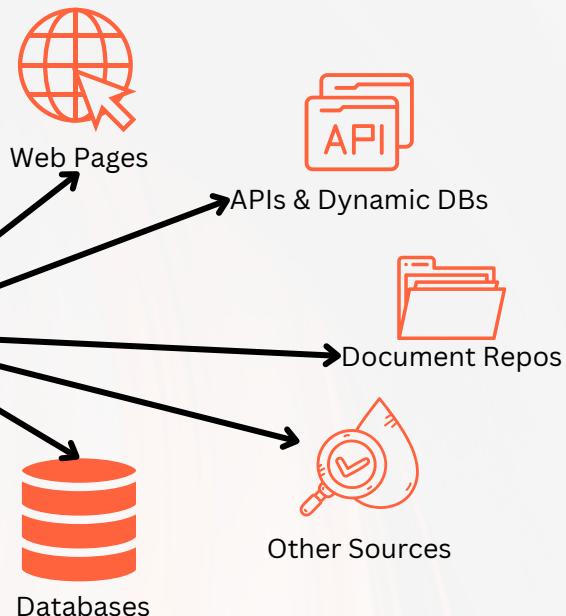
The Retriever of an RAG system can have access to external sources of information. Therefore, the LLM is not limited to its internal knowledge. The external sources can be proprietary documents and data or even the internet.

Without RAG With RAG



An LLM has knowledge only of the data it has been trained on
Also called **Parametric Memory** (information stored in the model parameters)

Retriever



Retriever searches and fetches information that the LLM has not necessarily been trained on. This adds to the LLM memory and is passed as context in the prompts. Also called **Non-Parametric Memory** (information available outside the model parameters)

- Expandable to all sources
- Easier to update/maintain
- Much cheaper than retraining/fine-tuning

The effort lies in creation of the knowledge base

Confidence in Responses

With the context (extra information that is retrieved) made available to the LLM, the confidence in LLM responses is increased.



Context Awareness

Added information assists LLMs in generating responses that are accurate and contextually appropriate



Source Citation

Access to sources of information improves the transparency of the LLM responses



Reduced Hallucinations

RAG enabled LLM systems are observed to be less prone to hallucinations than the ones without RAG

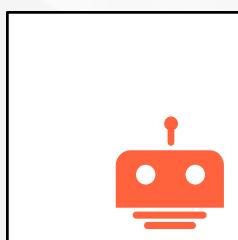
RAG Use Cases

The development of RAG technique is rooted in use cases that were limited by the inherent weaknesses of the LLMs. As of today some commercial applications of RAG are in -



Document Question Answering Systems

By providing access to proprietary enterprise document to an LLM, the responses are limited to what is provided within them. A retriever can search for the most relevant documents and provide the information to the LLM. Check out [this blog](#) for an example



Conversational agents

LLMs can be customised to product/service manuals, domain knowledge, guidelines, etc. using RAG. The agent can also route users to more specialised agents depending on their query. [SearchUnify has an LLM+RAG powered conversational agent](#) for their users.



Real-time Event Commentary

Imagine an event like a sports or a new event. A retriever can connect to real-time updates/data via APIs and pass this information to the LLM to create a virtual commentator. These can further be augmented with Text To Speech models. [IBM leveraged the technology for commentary during the 2023 US Open](#)



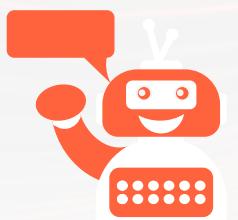
Content Generation

The widest use of LLMs has probably been in content generation. Using RAG, the generation can be personalised to readers, incorporate real-time trends and be contextually appropriate. [Yarnit is an AI based content marketing platform that uses RAG for multiple tasks](#).



Personalised Recommendation

Recommendation engines have been a game changer in the digital economy. LLMs are capable of powering the next evolution in content recommendations. Check out [Aman's blog](#) on the utility of LLMs in recommendation systems.

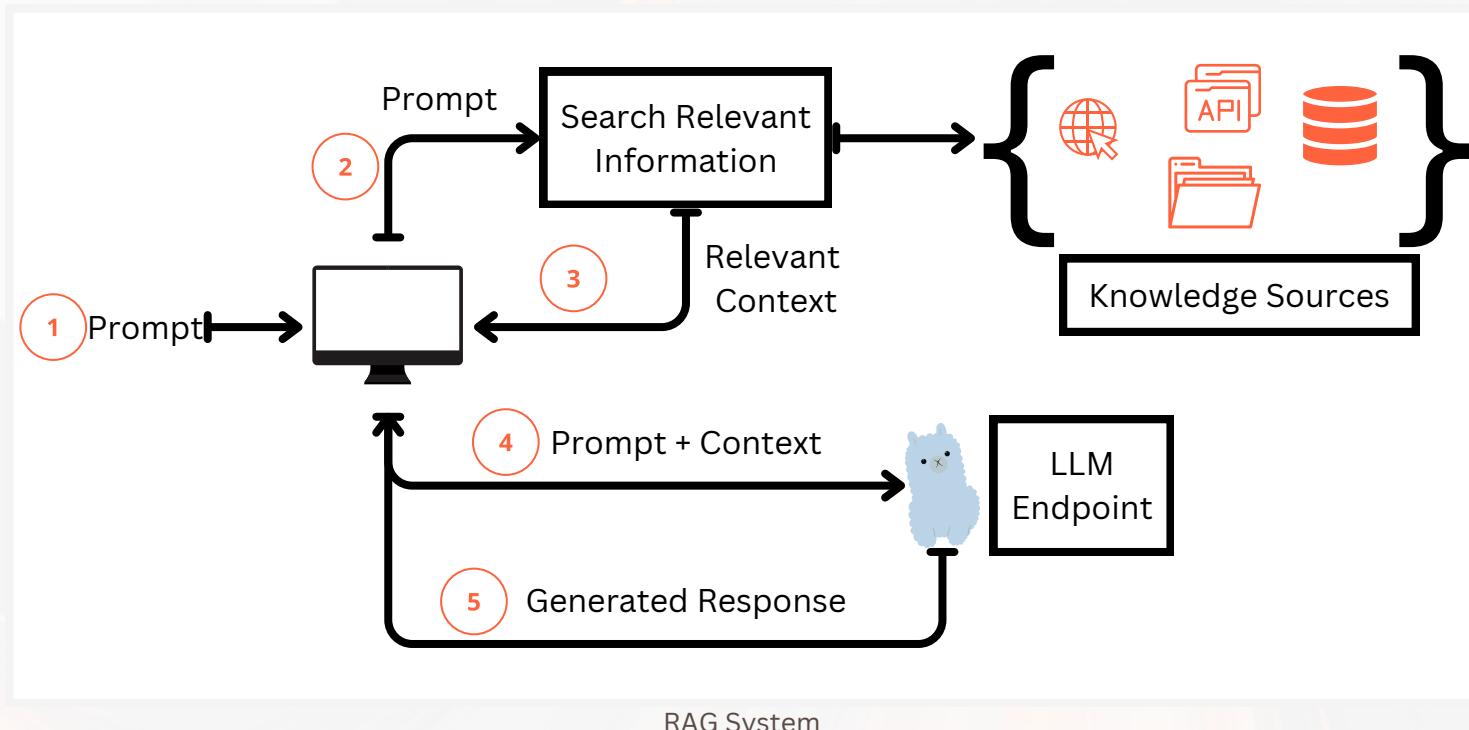


Virtual Assistants

Virtual personal assistants like Siri, Alexa and others are in plans to use LLMs to enhance the experience. Coupled with more context on user behaviour, these assistants can become highly personalised.

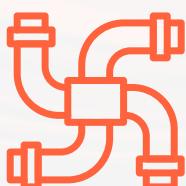
RAG Architecture

Let's revisit the five high level steps of an RAG enabled system



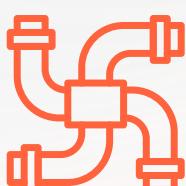
- 1 User writes a prompt or a query that is passed to an orchestrator
- 2 Orchestrator sends a search query to the retriever
- 3 Retriever fetches the relevant information from the knowledge sources and sends back
- 4 Orchestrator augments the prompt with the context and sends to the LLM
- 5 LLM responds with the generated text which is displayed to the user via the orchestrator

Two pipelines become important in setting up the RAG system. The first one being setting up the knowledge sources for efficient search and retrieval and the second one being the five steps of the generation.



Indexing Pipeline

Data for the knowledge is ingested from the source and indexed. This involves steps like splitting, creation of embeddings and storage of data.

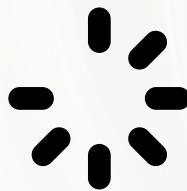


RAG Pipeline

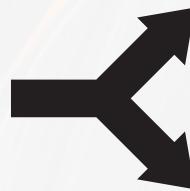
This involves the actual RAG process which takes the user query at run time and retrieves the relevant data from the index, then passes that to the model

Indexing Pipeline

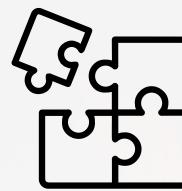
The indexing pipeline sets up the knowledge source for the RAG system. It is generally considered an offline process. However, information can also be fetched in real time. It involves four primary steps.



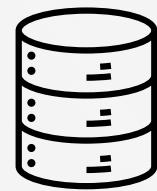
Loading



Splitting



Embedding



Storing

This step involves extracting information from different knowledge sources and loading them into documents.

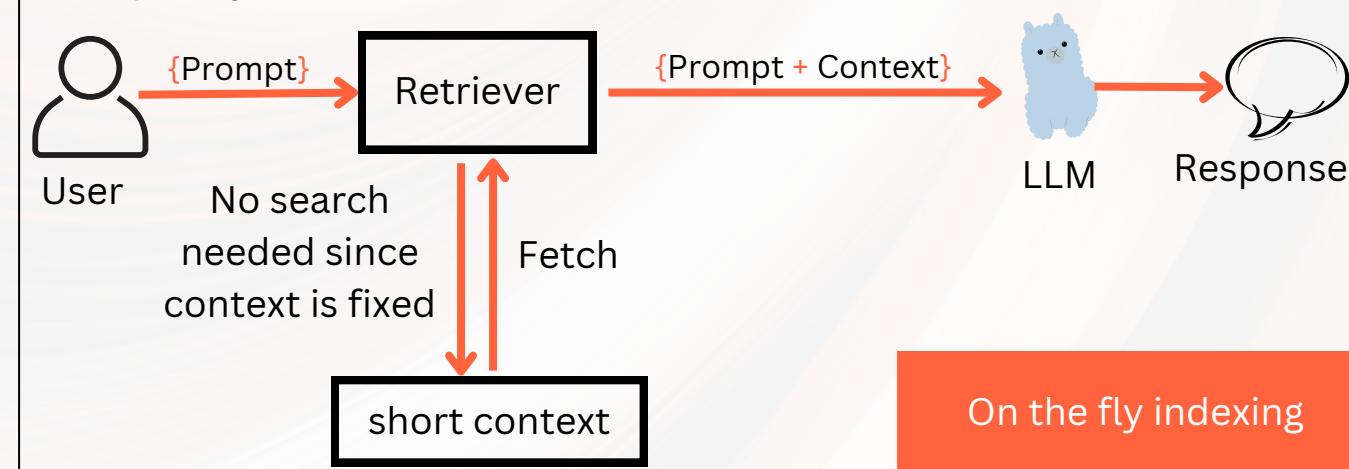
This step involves splitting documents into smaller manageable chunks. Smaller chunks are easier to search and to use in LLM context windows.

This step involves converting text documents into numerical vectors. ML models are mathematical models and therefore require numerical data.

This step involves storing the embeddings vectors. Vectors are typically stored in Vector Databases which are best suited for searching.

Offline Indexing pipelines are typically used when a knowledge base with large amount of data is being built for repeated usage e.g. a number of enterprise documents, manuals etc.

In cases where only a fixed small amount of one time data is required e.g. a 300 word blog, there is no need for storing the data. The blog text can either be directly passed in the LLM context window or a temporary vector index can be created.



Loading Data

As we've been discussing, the utility of RAG is to access data for all sorts of sources. These sources can be -

- Websites & HTML pages
- Documents like word, pdf etc.
- Code in python, java etc.
- Data in json, csv etc.
- APIs
- File Directories
- Databases
- And many more

The first step is to extract the information present in these source locations.

This is a good time to introduce two popular frameworks that are being used to develop LLM powered applications.



LangChain

Use cases: Good for applications that need enhanced AI capabilities, like language understanding tasks and more sophisticated text generation

Features: Stands out for its versatility and adaptability in building robust applications with LLMs

Agents: Makes creating agents using large language models simple through their agents API



Llamaindex

Use cases: Good for tasks that require text search and retrieval, like information retrieval or content discovery

Features: Excels in data indexing and language model enhancement

Connectors: Provides connectors to access data from databases, external APIs, or other datasets

Both frameworks are rapidly evolving and adding new capabilities every week. It's not an either/or situation and you can use both together (or neither).

Example : Loading a YouTube Video Transcript using LangChain Loaders

Let's begin by sourcing the transcript from this video -

“DALL·E 2 Explained” by OpenAI

(<https://www.youtube.com/watch?v=qTgPSKKjfVg>)

Below is the code using YoutubeLoader from langchain.document_loaders

```
● ● ●

from langchain.document_loaders import YoutubeLoader

loader = YoutubeLoader.from_youtube_url(
    "https://www.youtube.com/watch?v=qTgPSKKjfVg",
    add_video_info=True
)

loader.load()
```

LangChain Document Loader : YoutubeLoader

Loader object

[Document(page_content="Have you ever seen a polar bear playing bass? Or a robot painted like a Picasso? Didn't think so. DALL-E 2 is

....

....

....humans\and clever systems can work together to make new things – amplifying our creative potential.", metadata={'source': 'qTgPSKKjfVg', 'title': 'DALL·E 2 Explained', 'description': 'Unknown', 'view_count': 853564, 'thumbnail_url': 'https://i.ytimg.com/vi/qTgPSKKjfVg/hq720.jpg', 'publish_date': '2022-04-06 00:00:00', 'length': 167, 'author': 'OpenAI'})]

The Document object contains the **page_content** which is the transcript extracted from the youtube video as well as the **metadata** description

Example : Loading a Webpage Text using LlamaIndex Reader

This is a blog published on Medium -

What is a fine-tuned LLM?

(<https://medium.com/mlearning-ai/what-is-a-fine-tuned-llm-67bf0b5df081>)

Below is the code using SimpleWebPageReader from llama_hub

```
from llama_hub.web.simple_web.base import SimpleWebPageReader

loader = SimpleWebPageReader()

docs = loader.load_data(
    urls=["https://medium.com/mlearning-ai/\
what-is-a-fine-tuned-llm-67bf0b5df081"]
)
```

LlamaIndex LlamaHub Web Page Reader

Loader object

```
[Document(id_='17761da4-6a3a-4ce5-8590-c65ee446788f',
embedding=None, metadata={}, excluded_embed_metadata_keys=[],
excluded_llm_metadata_keys=[], relationships={},
hash='6471b3ffe4d3abb1aba2ca99d1d0448e2c3cbd157ddca256fab9fa363e0
9ed85', text='<!doctype html><html lang="en"><head><title data-
rh="true">What is a fine-tuned LLM?. Fine-tuning large language models...
| by Abhinav Kimothi |
...
</body></html>', start_char_idx=None, end_char_idx=None,
text_template='{metadata_str}\n\n{content}', metadata_template='{key}:
{value}', metadata_seperator='\n')]
```

The LlamaIndex Document object contains more attributes than a LangChain Document. Apart from **text** and **metadata**, it also has **id**, **templates** and other customizations available

Both LangChain and LlamaIndex offer loader integrations with more than a hundred data sources and the list keeps on growing

LangChain Document Loaders

The screenshot shows the LangChain documentation page for 'Document loaders'. It lists several pre-built loader components:

- Diffbot**: Unlike traditional web scraping tools, Diffbot doesn't require any...
- Discord**: Discord is a VoIP and instant messaging social platform. Users ha...
- Docugami**: This notebook covers how to load documents from Docugami. It...
- Docusaurus**: Docusaurus is a static-site generator which provides out-of-the-
- Dropbox**: Dropbox is a file hosting service that brings everything tradition...
- DuckDB**: DuckDB is an in-process SQL OLAP database management system...
- Email**: embass is a fully managed NLP API service that offers features li...
- Embas**
- EPub**: EPUB is an e-book file format that uses the "epub" file extension...
- Etherscan**: Etherscan is the leading blockchain explorer, search, API and anal...

LangChain provides integrations with a variety of sources

LlamaHub Data Loaders

The screenshot shows the LlamaHub 'Data Loaders' page. It displays a grid of data source cards:

Category	Source	Author	Last Updated
database	database	kevinpig	2 months ago
	database	jerryliu	2 months ago
motion	motion	kevinpig	5 months ago
	motion	jerryliu	2 months ago
google_docs	google_docs	kevinpig	2 months ago
	google_docs	theveseshang	15 days ago
slack	slack	kevinpig	1 month ago
	slack	jerryliu	1 month ago
mongo	mongo	kevinpig	24 days ago
	mongo	jerryliu	4 months ago
weaviate	weaviate	kevinpig	4 months ago
	weaviate	jerryliu	6 months ago
discord	discord	kevinpig	2 months ago
	discord	jerryliu	2 months ago
papers/arxiv	papers/arxiv	theveseshang	2 months ago
	papers/arxiv	theveseshang	5 months ago
file/docx	file/docx	kevinpig	2 months ago
	file/docx	theveseshang	13 days ago
twitter	twitter	kevinpig	5 months ago
	twitter	theveseshang	12 days ago
papers/pubmed	papers/pubmed	theveseshang	2 months ago
	papers/pubmed	theveseshang	5 months ago
file/audio	file/audio	kevinpig	13 days ago
	file/audio	theveseshang	4 months ago
file/pptx	file/pptx	kevinpig	5 months ago
	file/pptx	theveseshang	5 months ago
file/epub	file/epub	kevinpig	5 months ago
	file/epub	theveseshang	7 days ago
make_com	make_com	kevinpig	2 months ago
	make_com	jerryliu	2 months ago
Talos	Talos	kevinpig	6 months ago
	Talos	jerryliu	6 months ago

LlamaIndex provides data loaders via LlamaHub

These Document Loaders are particularly helpful in quickly making connections and accessing information. For specific sources, custom loaders can also be developed.

It is worthwhile exploring documentation for both

LlamaIndex: <https://docs.llamaindex.ai/en/stable/>

LangChain: https://python.langchain.com/docs/get_started/introduction

- *Loading documents from a list of sources may turn out to be a complicated process. Make sure to plan for all the sources and loaders in advance.*
- *More often than naught, transformations/clean-ups to the loaded data will be required like removing duplicate content, html parsing, etc. LangChain also provides a variety of document transformers*

Document Splitting

Once the data is loaded, the next step in the indexing pipeline is splitting the documents into manageable chunks. The question arises around the need of this step. Why is splitting of documents necessary. There are two reasons for that -

Ease of Search

Large chunks of data are harder to search over. Splitting data into smaller chunks therefore helps in better indexation.



Context Window Size

LLMs allow only a finite number of tokens in prompts and completions. The context therefore cannot be larger than what the context window permits.

Chunking Strategies

While splitting documents into chunks might sound a simple concept, there are certain best practices that researchers have discovered. There are a few considerations that may influence the overall chunking strategy.

1 Nature of Content

Consider whether you are working with lengthy documents, such as articles or books, or shorter content like tweets or instant messages. The chosen model for your goal and, consequently, the appropriate chunking strategy depend on your response.

2 Embedding Model being Used

We will discuss embeddings in detail in the next section but the choice of embedding model also dictates the chunking strategy. Some models perform better with chunks of specific length

3 Expected Length and Complexity of User Queries

Determine whether the content will be short and specific or long and complex. This factor will influence the approach to chunking the content, ensuring a closer correlation between the embedded query and the embedded chunks

4 Application Specific Requirements

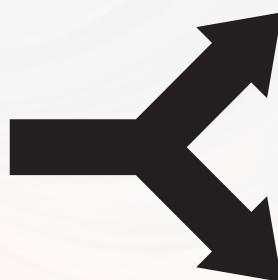
The application use case, such as semantic search, question answering, summarization, or other purposes will also determine how text should be chunked. If the results need to be input into another language model with a token limit, it is crucial to factor this into your decision-making process.

Chunking Methods

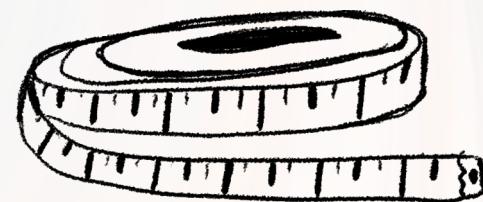
Depending on the aforementioned considerations, a number of **text splitters** are available. At a broad level, text splitters operate in the following manner:

- Divide the text into compact, **semantically meaningful units**, often sentences.
- Merge these smaller units into larger chunks until a specific size is achieved, measured by a **length function**.
- Upon reaching the predetermined size, treat that chunk as an independent segment of text. Thereafter, start creating a new text chunk with **some degree of overlap** to maintain contextual continuity between chunks.

Two areas to focus on, therefore are -



How the text is split?



How the chunk size is measured?

A very common approach is where we **pre-determine** the size of the text chunks.

Additionally, we can specify the **overlap between chunks** (Remember, overlap is preferred to maintain contextual continuity between chunks).

This approach is simple and cheap and is, therefore, widely used. Let's look at some examples -

Split by Character

In this approach, the text is split based on a character and the chunk size is measured by the number of characters.

Example text : alice_in_wonderland.txt (the book in .txt format)

using **LangChain's CharacterTextSplitter**

```
● ● ●

with open('../Data/alice_in_wonderland.txt') as f:
    AliceInWonderland = f.read()

from langchain.text_splitter import CharacterTextSplitter
text_splitter = CharacterTextSplitter(
    separator = "\n\n",
    chunk_size = 2000,
    chunk_overlap = 100,
    length_function = len,
    is_separator_regex = False,
)

texts = text_splitter.create_documents([AliceInWonderland])
print(texts[0])
print(texts[1])
```

texts[0]

"TITLE: Alice's Adventures in Wonderland\nAUTHOR: Lewis Carroll\n\n CHAPTER I \n(Down the Rabbit-Hole)\n\n Alice was beginning to get very tired of sitting by her sister\non the bank, and of having nothing to do: once or twice she had\npeeped into the book her sister was reading, but it had no\npictures or conversations in it, and what is the use of a book,' thought Alice `without pictures or conversation?'\n\n So she was considering in her own mind (as well as she could,\nfor the hot day made her feel very sleepy and stupid), whether\nthe pleasure of making a daisy-chain would be worth the trouble\nof getting up and picking the daisies, when suddenly a White\nRabbit with pink eyes ran close by her.\n\n There was nothing so VERY remarkable in that; nor did Alice\nthink it so VERY much out of the way to hear the Rabbit say to\nitself, `Oh dear! Oh dear! I shall be late!' (when she thought\nit over afterwards, it occurred to her that she ought to have\nwondered at this, but at the time it all seemed quite natural);\nbut when the Rabbit actually TOOK A WATCH OUT OF ITS WAISTCOAT-\nPOCKET, and looked at it, and then hurried on, Alice started to\nher feet, for it flashed across her mind that she had never\nbefore seen a rabbit with either a waistcoat-pocket, or a watch to\ntake out of it, and burning with curiosity, she ran across\nthe field after it, and fortunately was just in time to see it pop\ndown a large rabbit-hole under the hedge.\n\n In another moment down went Alice after it, never once\nconsidering how in the world she was to get out again.\n\n The rabbit-hole went straight on like a tunnel for some way,\nand then dipped suddenly down, so suddenly that Alice had not\na moment to think about stopping herself before she found herself\nfalling down a very deep well."

Overlap

texts[1]

"In another moment down went Alice after it, never once\nconsidering how in the world she was to get out again.\n\n The rabbit-hole went straight on like a tunnel for some way,\nand then dipped suddenly down, so suddenly that Alice had not\na moment to think about stopping herself before she found herself\nfalling down a very deep well.\n\n Either the well was very deep, or she fell very\nslowly, for she\nhad plenty of time as she went down to look about her and to\nwonder what was going to happen next. First, she tried to look\ndown and make out what she was coming to, but it was too dark to\nsee anything; then she looked at the sides of the well, and\nnoticed that they were filled with cupboards and book-shelves;\nhere and there she saw maps and pictures hung upon pegs. She\ntook down a jar from one of the shelves as she passed; it was\nlabelled `ORANGE MARMALADE'; but to her great disappointment\nit was empty: she did not like to drop the jar for fear of killing\nsomebody, so managed to put it into one of the cupboards as she\nfell past it."

Let's find out how many chunks were created

```
● ● ●
print(f"Total Number of Chunks Created => {len(texts)}")
print(f"Length of the First Chunk is => {len(texts[0].page_content)} characters")
print(f"Length of the Last Chunk is => {len(texts[-1].page_content)} characters")
```

Total Number of Chunks Created => 93

Length of the First Chunk is => 1777 characters

Length of the Last Chunk is => 816 characters

Recursive Split by Character

A subtle variation to splitting by character is Recursive Split. The only difference is that instead of a single character used for splitting, this technique **uses a list of characters** and tries to split hierarchically till the chunk sizes are small enough. This technique is generally recommended for generic text.

Example text : AK_BusyPersonIntroLLM.txt

(Transcript of a YouTube video by Andrej Karpathy titled [1hr Talk] Intro to Large Language Models - https://www.youtube.com/watch?v=zjkBMFhNj_g&t=9s)

using **LangChain's RecursiveCharacterTextSplitter**

This is a generic text that is not formatted. Let's compare the two strategies.

with CharacterTextSplitter

```
● ● ●
with open('../Data/AK_BusyPersonIntroLLM.txt') as f:
    IntroToLLM = f.read()

from langchain.text_splitter import CharacterTextSplitter
text_splitter = CharacterTextSplitter(
    separator = "\n\n",
    chunk_size = 2000,
    chunk_overlap = 400,
    length_function = len,
    is_separator_regex = False,
)

texts = text_splitter.create_documents([IntroToLLM])

print(f"Total Number of Chunks Created => {len(texts)}")
print(f"Length of the First Chunk is => {len(texts[0].page_content)} characters")
print(f"Length of the Last Chunk is => {len(texts[-1].page_content)} characters")
```

Total Number of Chunks Created => 1

Length of the First Chunk is => 64383 characters

Length of the Last Chunk is => 64383 characters

Text splitter fails to convert the text into chunks since there are no '\n\n' character present in the raw transcript

with RecursiveCharacterTextSplitter

```
● ● ●  
with open('../Data/AK_BusyPersonIntroLLM.txt') as f:  
    IntroToLLM = f.read()  
  
from langchain.text_splitter import RecursiveCharacterTextSplitter  
  
text_splitter = RecursiveCharacterTextSplitter(  
    chunk_size = 2000,  
    chunk_overlap = 400,  
    length_function = len,  
    is_separator_regex = False,  
)  
  
texts = text_splitter.create_documents([IntroToLLM])  
  
print(f"Total Number of Chunks Created => {len(texts)}")  
  
print(f"Length of the First Chunk is => {len(texts[0].page_content)} characters")  
  
print(f"Length of the Last Chunk is => {len(texts[-1].page_content)} characters")
```

Total Number of Chunks Created => 40
Length of the First Chunk is => 1998 characters
Length of the Last Chunk is => 1967 characters

*Recursive text splitter performs well in dealing
with generic text*

Split by Tokens

For those well versed with Large Language Models, tokens is not a new concept. All LLMs have a token limit in their respective context windows which we cannot exceed. It is therefore a good idea to count the tokens while creating chunks. All LLMs also have their tokenizers.

Tiktoken Tokenizer



Tiktoken tokenizer has been created by OpenAI for their family of models. Using this strategy, the split still happens based on the character. However, the length of the chunk is determined by the number of tokens.

Example text : AK_BusyPersonIntroLLM.txt

(Transcript of a YouTube video by Andrej Karpathy titled [1hr Talk] Intro to Large Language Models - https://www.youtube.com/watch?v=zjkBMFhNj_g&t=9s)

using LangChain's TokenTextSplitter

```
● ● ●  
with open('../Data/AK_BusyPersonIntroLLM.txt') as f:  
    IntroToLLM = f.read()  
  
from langchain.text_splitter import TokenTextSplitter  
import tiktoken  
  
text_splitter = TokenTextSplitter(  
    chunk_size = 1000,  
    chunk_overlap = 20,  
    length_function = len,  
)  
  
texts = text_splitter.create_documents([IntroToLLM])  
  
encoding = tiktoken.get_encoding("cl100k_base")  
  
print(f"Total Number of Chunks Created => {len(texts)}")  
print(f"Total Number of Tokens in the document => {len(encoding.encode(IntroToLLM))} tokens")  
print(f"Length of the First Chunk is => {len(encoding.encode(texts[0].page_content))} tokens")  
print(f"Length of the Last Chunk is => {len(encoding.encode(texts[-1].page_content))} tokens")
```

Total Number of Chunks Created => 14
Total Number of Tokens in the document => 12865 tokens
Length of the First Chunk is => 1014 tokens
Length of the Last Chunk is => 1014 tokens

Tokenizers are helpful in creating chunks that sit well in the context window of an LLM

Hugging Face Tokenizer



Hugging Face has become the go-to platform for anyone building apps using LLMs or even other models. All models available via Hugging Face are also accompanied by their tokenizers.

Example text : AK_BusyPersonIntroLLM.txt

(Transcript of a YouTube video by Andrej Karpathy titled [1hr Talk] Intro to Large Language Models - https://www.youtube.com/watch?v=zjkBMFhNj_g&t=9s)

using **Transformers** and **LangChain's RecursiveCharacterTextSplitter**

Example tokenizer : GPT2TokenizerFast

```
● ● ●

with open('../Data/AK_BusyPersonIntroLLM.txt') as f:
    IntroToLLM = f.read()

from transformers import GPT2TokenizerFast
from langchain.text_splitter import RecursiveCharacterTextSplitter

tokenizer = GPT2TokenizerFast.from_pretrained("gpt2")

text_splitter = RecursiveCharacterTextSplitter.from_huggingface_tokenizer(
    tokenizer, chunk_size=100, chunk_overlap=0
)

texts = text_splitter.split_text(IntroToLLM)

print(texts[0])
print(texts[1])
```

texts[0]

"hi everyone so recently I gave a 30-minute talk on large language models just kind of like an intro talk um unfortunately that talk was not recorded but a lot of people came to me after the talk and they told me that uh they liked the talk so I would just I thought I would just re-record it and basically put it up on YouTube so here we go the busy person's intro to large language models director Scott okay so let's begin first of all what is a large language model

No Overlap as specified

texts[1]

really well a large language model is just two files right um there be two files in this hypothetical directory so for example work with the specific example of the Llama 270b model this is a large language model released by meta AI and this is basically the Llama series of language models the second iteration of it and this is the 70 billion parameter model of uh of this series so there's multiple models uh belonging to the Llama 2 Series uh 7 billion um 13 billion 34 billion and 70 billion is the the

Do take a look at Hugging Face documents on Tokenizers



Hugging Face

https://huggingface.co/docs/transformers/tokenizer_summary

Other Tokenizer

Other libraries like Spacy, NLTK and SentenceTransformers also provide splitters



```
from langchain.text_splitter import NLTKTextSplitter
text_splitter = NLTKTextSplitter(chunk_size=1000)
texts = text_splitter.split_text(IntroToLLM)
print(texts[0])
```



spaCy



```
from langchain.text_splitter import SpacyTextSplitter
text_splitter = SpacyTextSplitter(chunk_size=1000)
texts = text_splitter.split_text(IntroToLLM)
print(texts[0])
```

Specialized Chunking

Chunking often aims to keep text with common context together. With this in mind, we might want to specifically honour the structure of the document itself for example **HTML**, **Markdown**, **Latex** or even **code**.

Example : <https://medium.com/p/29a7e8610843>

Example HTML : “Context is Key: The Significance of RAG in Language Models”

(A blog on Medium - <https://medium.com/p/29a7e8610843>)

using **LangChain’s HTMLHeaderTextSplitter & RecursiveCharacterTextSplitter**

```
from langchain.text_splitter import HTMLHeaderTextSplitter
from langchain.text_splitter import RecursiveCharacterTextSplitter
url = "https://medium.com/p/29a7e8610843"

headers_to_split_on = [
    ("h1", "Header 1"),
    ("h2", "Header 2"),
    ("h3", "Header 3"),
    ("h4", "Header 4"),
]

html_splitter = HTMLHeaderTextSplitter(headers_to_split_on=headers_to_split_on)

# for local file use html_splitter.split_text_from_file(<path_to_file>
html_header_splits = html_splitter.split_text_from_url(url)

chunk_size = 5000
chunk_overlap = 300
text_splitter = RecursiveCharacterTextSplitter(
    chunk_size=chunk_size, chunk_overlap=chunk_overlap
)

# Split
splits = text_splitter.split_documents(html_header_splits)
len(splits)
```



LangChain [All LangChain Splitters](#)

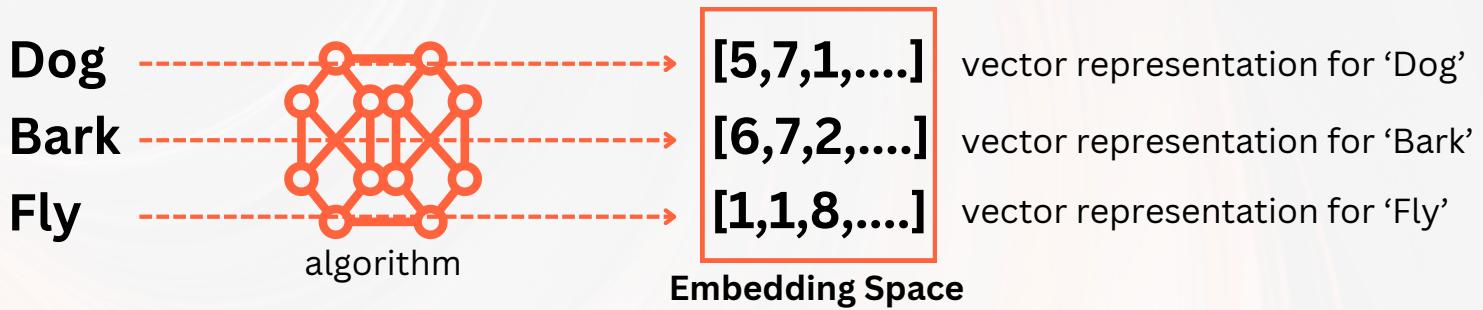


Things to Keep in Mind

- Ensure data quality by preprocessing it before determining the optimal chunk size. Examples include removing HTML tags or eliminating specific elements that contribute noise, particularly when data is sourced from the web.
- Consider factors such as content nature (e.g., short messages or lengthy documents), embedding model characteristics, and capabilities like token limits in choosing chunk sizes. Aim for a balance between preserving context and maintaining accuracy.
- Test different chunk sizes. Create embeddings for the chosen chunk sizes and store them in your index or indices. Run a series of queries to evaluate quality and compare the performance of different chunk sizes.

Embeddings

All Machine Learning/AI models work with numerical data. Before the performance of any operation all text/image/audio/video data has to be transformed into a numerical representation. **Embeddings are vector representations of data that capture meaningful relationships between entities.** As a general definition, embeddings are data that has been transformed into n-dimensional matrices for use in deep learning computations. A word embedding is a vector representation of words.



The process of embedding **transforms** data (like text) into vectors, **compresses** the input information resulting in an **embedding space specific to the training data**

While we keep our discussion around embeddings limited to RAG application and how to create embeddings for our data, a great resource to find more about embeddings is this book by Vicki Boykis [[What are embeddings](#)]



What are embeddings

Vicki Boykis

The good news for anyone building RAG Applications is that embeddings once created can also generalize to other tasks and domains through **transfer learning** – the ability to switch contexts – which is one of the reasons embeddings have exploded in popularity across machine learning applications

Popular Embedding Models

word2vec

Google's Word2Vec is one of the most popular pre-trained word embeddings. The official paper - <https://arxiv.org/pdf/1301.3781.pdf>

GLOVE

The 'Global Vectors' model is so termed because it captures statistics directly at a global level. The official paper - <https://nlp.stanford.edu/pubs/glove.pdf>

fastText

Facebook's AI research, fastText builds embeddings composed of characters instead of words. The official paper - <https://arxiv.org/pdf/1607.04606.pdf>

Elmo

Embeddings from Language Models, are learnt from the internal state of a bidirectional LSTM. The official paper - <https://arxiv.org/pdf/1802.05365.pdf>

BERT

Bidirectional Encoder Representations from Transformers is a transformer bases approach. The official paper - <https://arxiv.org/pdf/1810.04805.pdf>

ada v2 by OpenAI

used by GPT series of models

textembedding-gecko

by Google's



Other Open Source Embeddings

Checkout [MTEB leaderboard](#) at



Hugging Face



Rank	Model	Model Size (GB)	Embedding Dimensions	Sequence Length	Average (56 datasets)	Classification Average (22 datasets)	Clustering Average (11 datasets)	Pair Classification Average (3 datasets)	Ranking Average (4 datasets)	Retrieval Average (15 datasets)	STS Average (10 datasets)	Sum Average datasets
1	UAE-Large-XL	1.34	1024	512	64.64	75.58	46.73	87.25	59.88	54.66	84.54	32.03
2	xoxoxo-01-instant-t		1024	4996	64.49	74.79	47.4	86.57	59.74	55.58	82.93	38.97
3	Cobain-embed-english-v3.0		1024	512	64.47	76.49	47.43	85.84	58.01	55	82.62	38.18
4	beto-large-en-v1.5	1.34	1024	512	64.23	75.97	46.08	87.12	60.03	54.29	83.11	31.61

How to Choose Embeddings?

Ever since the release of ChatGPT and the advent of the aptly described LLM Wars, there has also been a mad rush in developing embeddings models. There are many evolving standards of evaluating LLMs and embeddings alike.

When building RAG powered LLM apps, there is no right answer to “Which embeddings model to use?”. However, you may notice particular embeddings working better for specific use cases (like summarization, text generations, classification etc.)

	Models	Use Cases
Text similarity: Captures semantic similarity between pieces of text.	text-similarity-{ada, babbage, curie, davinci}-001	Clustering, regression, anomaly detection, visualization
Text search: Semantic information retrieval over documents.	text-search-{ada, babbage, curie, davinci}-{query, doc}-001	Search, context relevance, information retrieval
Code search: Find relevant code with a query in natural language.	code-search-{ada, babbage}-{code, text}-001	Code search and relevance

OpenAI used to recommend different embeddings models for different use cases. However, now they recommend **ada v2** for all tasks.

Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)	Pair Classification Average (3 datasets)	Reranking Average (4 datasets)	Retrieval Average (15 datasets)	STS Average (10 datasets)	Summarization Average (1 dataset)
-----------------------	--------------------------------------	----------------------------------	--	--------------------------------	---------------------------------	---------------------------	-----------------------------------

MTEB Leaderboard at Hugging Face evaluates almost all available embedding models across seven use cases - *Classification*, *Clustering*, *Pair Classification*, *Reranking*, *Retrieval*, *Semantic Textual Similarity (STS)* and *Summarization*.

Another important consideration is **cost**. With OpenAI models you can incur significant costs if you are working with a lot of documents. The cost of open source models will depend on the implementation.

Creating Embeddings

Once you've chosen your embedding model, there are several ways of creating the embeddings. Sometimes, our friends, LlamaIndex and LangChain come in pretty handy to convert documents (*split into chunks*) into vector embeddings. Other times you can use the service from a provider directly or get the embeddings from HuggingFace

Example : OpenAI text-embedding-ada-002

using Embedding.create() function from openai library

```
● ● ●

with open('../Data/AK_BusyPersonIntroLLM.txt') as f:
    IntroToLLM = f.read()

from langchain.text_splitter import TokenTextSplitter

text_splitter = TokenTextSplitter(
    chunk_size = 1000,
    chunk_overlap = 20,
    length_function = len,
)

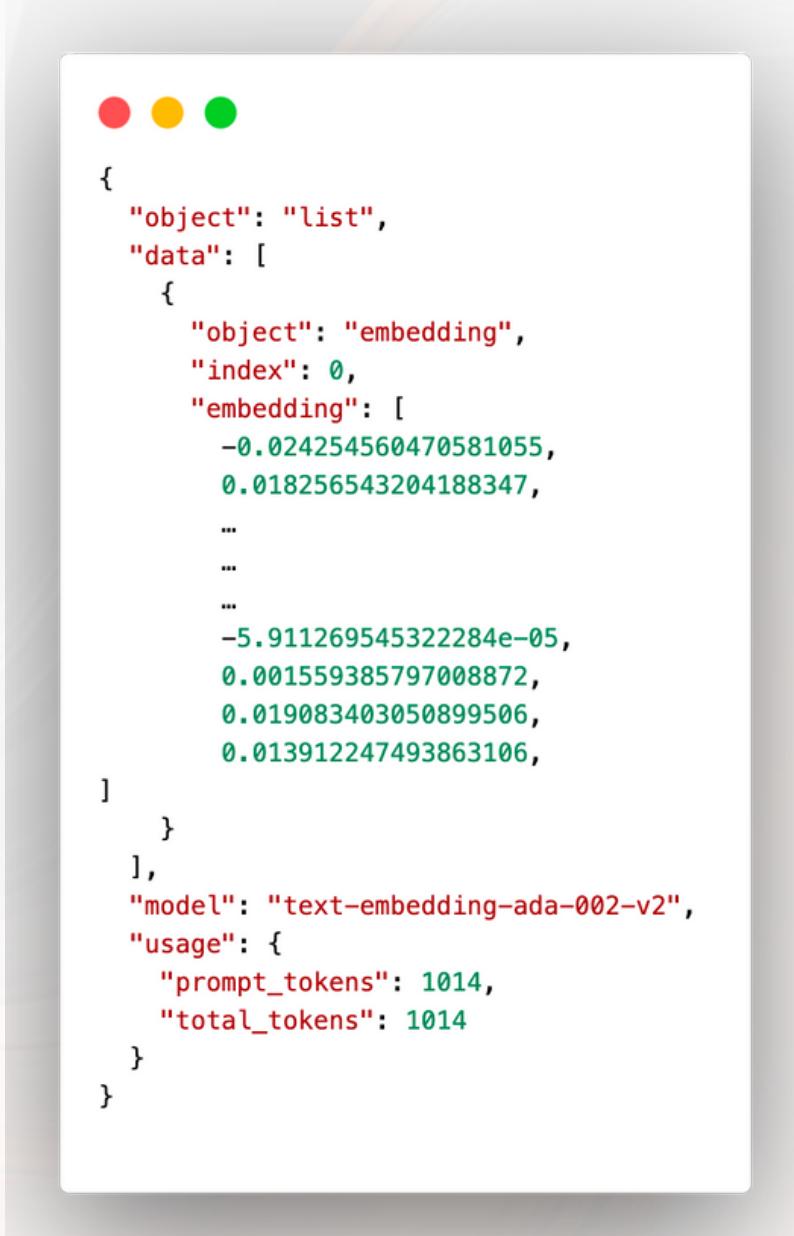
texts = text_splitter.create_documents([IntroToLLM])

import openai
openai.api_key='sk-#####eGeY#####RZj'

response = openai.Embedding.create(
    input=texts[0].page_content,
    model="text-embedding-ada-002"
)

print(response)
```

*You'll need an OpenAI apikey to create these embeddings
You can get one here - <https://platform.openai.com/api-keys>*

Example Response

```
{  
    "object": "list",  
    "data": [  
        {  
            "object": "embedding",  
            "index": 0,  
            "embedding": [  
                -0.024254560470581055,  
                0.018256543204188347,  
                ...  
                ...  
                ...  
                -5.911269545322284e-05,  
                0.001559385797008872,  
                0.019083403050899506,  
                0.013912247493863106,  
            ]  
        },  
        {  
            "model": "text-embedding-ada-002-v2",  
            "usage": {  
                "prompt_tokens": 1014,  
                "total_tokens": 1014  
            }  
        }  
    ]  
}
```

response.data[0].embedding will give the created embeddings that can be stored for retrieval

Cost

Model	Usage
ada v2	\$0.0001 / 1K tokens

In this example, 1014 tokens will cost about \$.0001. Recall that for this youtube transcript we got 14 chunks. So creating the embeddings for the entire transcript will cost about 0.14 cents. This may seem low, but when you scale up to thousands of documents being updated frequently, the cost can become a concern.

Example : msmarco-bert-base-dot-v5

using HuggingFaceEmbeddings from langchain.embeddings

```
from langchain.embeddings import HuggingFaceEmbeddings

embeddings = HuggingFaceEmbeddings(
    model_name='sentence-transformers/msmarco-bert-base-dot-v5'
)
text = texts[0].page_content

query_result = embeddings.embed_query(text)
```

Example : embed-english-light-v3.0

using CohereEmbeddings from langchain.embeddings

```
from langchain.embeddings import CohereEmbeddings

embeddings = CohereEmbeddings(
    model="embed-english-light-v3.0"
)

text = texts[0].page_content

query_result = embeddings.embed_query(text)
```



LangChain

All the available embeddings classes on LangChain



Storing

We are at the last step of creating the indexing pipeline. We have loaded and split the data, and created the embeddings. Now, for us to be able to use the information repeatedly, we need to store it so that it can be accessed on demand. For this we use a special kind of database called the **Vector Database**.

What is a Vector Database?

For those familiar with databases, **indexing** is a data structure technique that allows users to quickly retrieve data from a database. Vector databases specialise in indexing and storing embeddings for **fast retrieval** and **similarity search**.

A stripped down variant of a Vector Database is a Vector Index like FAISS (Facebook AI Similarity Search). It is this vector indexing that improves the search and retrieval of vector embeddings. Vector Databases augment the indexing with typical database features like data management, metadata storage, scalability, integrations, security etc.

In short, Vector Databases provide -

- Scalable Embedding Storage.
- Precise Similarity Search.
- Faster Search Algorithm.

Popular Vector Databases



Facebook AI Similarity search is a vector index released with a library in 2017



Pinecone is one of the most popular managed Vector DB for large scale



Weaviate

Weaviate is an open source vector database that stores both objects and vectors



Chromadb is also an open source vector database.

With the growth in demand for vector storage, it can be anticipated that all major database players will add the vector indexing capabilities to their offerings.

How to choose a Vector Database?

All vector databases offer the same basic capabilities. Your choice should be influenced by the nuance of your use case matching with the value proposition of the database.

A few things to consider -

- Balance search accuracy and query speed based on application needs. Prioritize accuracy for precision applications or speed for real-time systems.
- Weigh increased flexibility vs potential performance impacts. More customization can add overhead and slow systems down.
- Evaluate data durability and integrity requirements vs the need for fast query performance. Additional persistence safeguards can reduce speed.
- Assess tradeoffs between local storage speed and access vs cloud storage benefits like security, redundancy and scalability.
- Determine if tight integration control via direct libraries is required or if ease-of-use abstractions like APIs better suit your use case.
- Compare advanced algorithm optimizations, query features, and indexing vs how much complexity your use case necessitates vs needs for simplicity.
- Cost considerations - while you may incur regular cost in a fully managed solution, a self hosted one might prove costlier if not managed well

User Friendly for PoCs



Higher Performance



Customization



There are many more Vector DBs. For a comprehensive understanding of the pros and cons of each, this [blog](#) is highly recommended

Storing Embeddings in Vector DBs

To store the embeddings, LangChain and LlamalIndex can be used for quick prototyping. The more nuanced implementation will depend on the choice of the DB, use case, volume etc.

Example : FAISS from langchain.vectorstores

In this example, we complete our indexing pipeline for one document.

1. Loading our text file using **TextLoader**,
2. Splitting the text into chunks using **RecursiveCharacterTextSplitter**,
3. Creating embeddings using **OpenAIEmbeddings**
4. Storing the embeddings into **FAISS** vector index



```
from langchain.document_loaders import TextLoader
from langchain.embeddings.openai import OpenAIEmbeddings
from langchain.text_splitter import RecursiveCharacterTextSplitter
from langchain.vectorstores import FAISS

loader=TextLoader('../Data/AK_BusyPersonIntroLLM.txt')
document=loader.load()
text_splitter = RecursiveCharacterTextSplitter(chunk_size=10000,
                                              chunk_overlap=2000)
docs = text_splitter.split_documents(document)
num_emb=len(docs)
embeddings = OpenAIEmbeddings(openai_api_key=openai_api_key)
db = FAISS.from_documents(docs, embeddings)
```

You'll have to address the following dependencies.

1. Install openai, tiktoken and faiss-cpu or faiss-gpu
2. Get an OpenAI API key

Now that our knowledge base is ready, let's quickly see it in action. Let's perform a search on the FAISS index we've just created.

Similarity search

In the YouTube video, for which we have indexed the transcript, Andrej Karpathy talks about the idea of LLM as an operating system. Let's perform a search on this.

Query : What did Andrej say about LLM operating system?



```
query="What did Andrej say about LLM operating system?"  
docs = db.similarity_search(query)  
docs[0].page_content
```



"operating system and um basically this process
is coordinating a lot of resources be they memory
or computational tools for problem solving so let's
think through based on everything I've shown you
what an LLM might look like in a few years it can
read and generate text it has a lot more knowledge
than any single human about all the subjects it can browse
the internet or reference local files uh through
retrieval augmented generation it can use existing
software infrastructure like calculator python Etc it
can see and generate images and videos it can hear and
speak and generate music it can think for a long time using
a system too it can maybe self-improve in some narrow domains
that have a reward function available maybe it can be customized
and fine-tuned to many specific tasks maybe there's lots of
llm experts almost uh living in an App Store that can sort of
coordinate uh for problem solving and so I see a lot of equivalence
between this new llm OS operating system and operating"

We can see here that out of the entire text, we have been able to retrieve the specific chunk talking about the LLM OS. We'll look at it in detail again in the RAG pipeline

Example : Chroma from langchain.vectorstores

1. Loading our text file using **TextLoader**,
2. Splitting the text into chunks using **RecursiveCharacterTextSplitter**,
3. Creating embeddings using **all-MiniLM-L6-v2**
4. Storing the embeddings into **Chromadb**

```
● ● ●  
from langchain.document_loaders import TextLoader  
from langchain.embeddings.sentence_transformer import SentenceTransformerEmbeddings  
from langchain.text_splitter import CharacterTextSplitter  
from langchain.vectorstores import Chroma  
  
# load the document and split it into chunks  
loader = TextLoader('../Data/AK_BusyPersonIntroLLM.txt')  
documents = loader.load()  
  
# split it into chunks  
text_splitter = RecursiveCharacterTextSplitter(chunk_size=1000,  
                                              chunk_overlap=200)  
  
docs = text_splitter.split_documents(documents)  
  
# create the open-source embedding function  
embedding_function = SentenceTransformerEmbeddings(model_name="all-MiniLM-L6-v2")  
  
# load it into Chroma  
db = Chroma.from_documents(docs, embedding_function)  
  
# query it  
query = "What did Andrej say about LLM operating system?"  
docs = db.similarity_search(query)  
  
# print results  
print(docs[0].page_content)
```

```
● ● ●  
...  
llm trying to page relevant information in and out of its  
context window to perform your task um and so a lot of  
other I think connections also exist I think there's  
equivalence of um multi-threading multiprocessing speculative  
execution uh there's equivalent of in the random access memory  
in the context window there's equivalence of user space and  
kernel space and a lot of other equivalents to today's  
operating systems that I didn't fully cover but fundamentally  
the other reason that I really like this analogy of llms kind  
of becoming a bit of an operating system ecosystem is that  
there are also some equivalence I think between the current  
operating systems and the uh and what's emerging today so for  
example in the desktop operating system space we have a few  
proprietary operating systems like Windows and Mac OS but we  
also have this open source ecosystem of a large diversity of  
operating systems based on Linux in the same way here we have  
some proprietary operating systems like GPT  
...
```



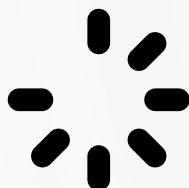
LangChain

All LangChain
VectorDB Integrations



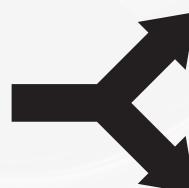
Indexing Pipeline Recap

We covered the indexing pipeline in its entirety. A quick recap -



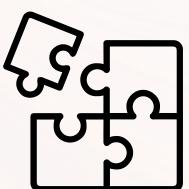
Loading

- A variety of data loaders from LangChain and LlamaIndex can be leveraged to load data from all sort of sources.
- Loading documents from a list of sources may turn out to be a complicated process. Make sure to plan for all the sources and loaders in advance.
- More often than naught, transformations/clean-ups to the loaded data will be required



Splitting

- Documents need to be split for ease of search and limitations of the llm context windows
- Chunking strategies are dependent on the use case, nature of content, embeddings, query length & complexity
- Chunking methods determine how the text is split and how the chunks are measured



Embedding

- Embeddings are vector representations of data that capture meaningful relationships between entities
- Some embeddings work better for some use cases



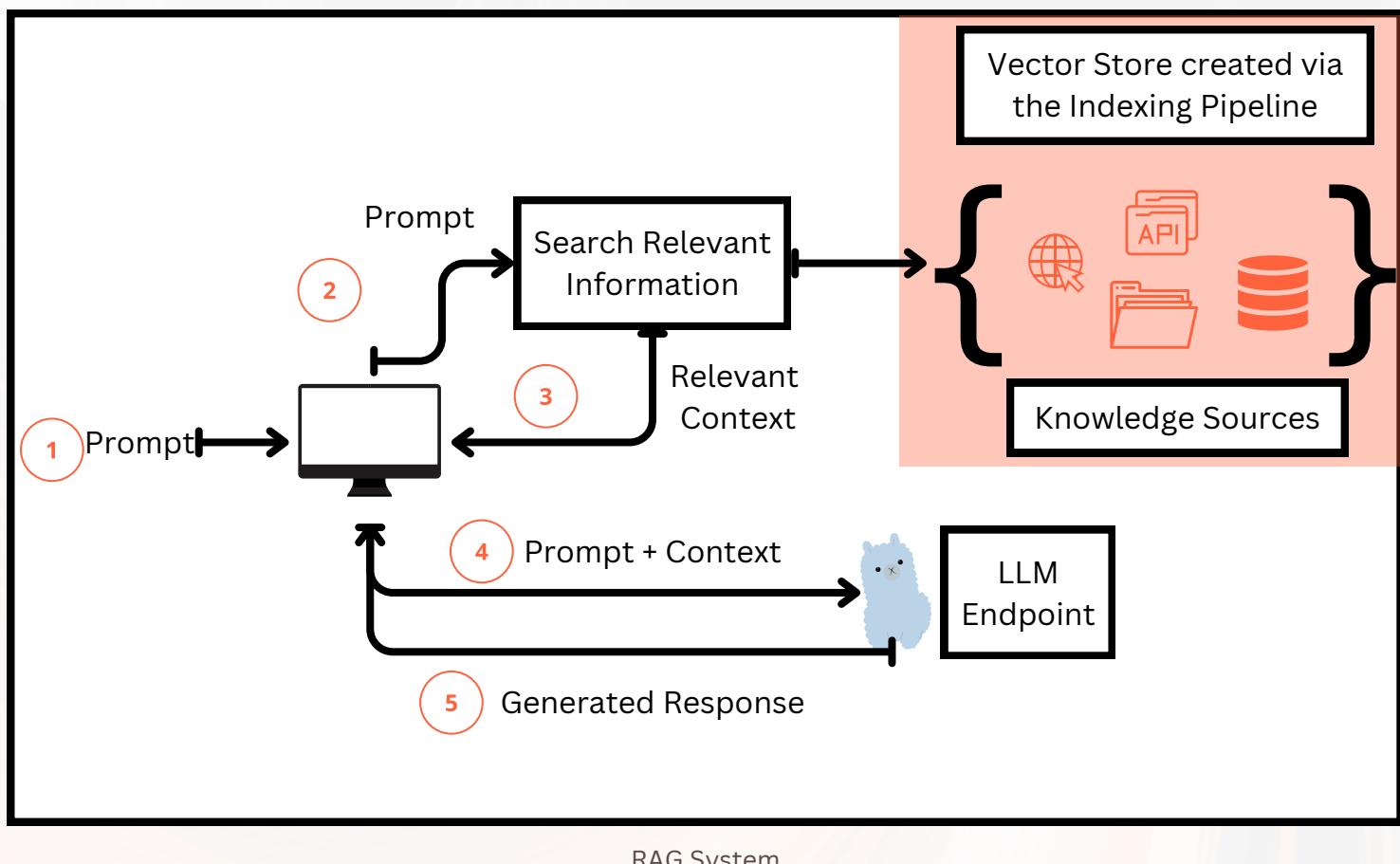
Storing

- Vector databases specialise in indexing and storing embeddings for fast retrieval and similarity search
- Different vector databases present different benefits and can be used in accordance with the use case

RAG Pipeline

Now that the knowledge base has been created in the indexing pipeline, the main generation or the **RAG pipeline** will have to be setup for receiving the input and generating the output.

Let's revisit our architecture diagram.



Generation Steps

- 1 User writes a prompt or a query that is passed to an orchestrator
- 2 Orchestrator sends a search query to the retriever
- 3 Retriever fetches the relevant information from the **knowledge sources** and returns
- 4 Orchestrator augments the prompt with the context and sends to the LLM
- 5 LLM responds with the generated text which is displayed to the user via the orchestrator

The knowledge sources highlighted above have been set up using the indexing pipeline. These sources can be served using “on-the-fly” indexing also

RAG Pipeline Steps

The three main steps in a RAG pipeline are



Search & Retrieval

This step involves searching for the context from the source (e.g. vector db)



Augmentation

This step involves adding the context to the prompt depending on the use case.



Generation

This step involves generating the final response from the large language model

An important consideration is how knowledge is stored and accessed. This has a bearing on the search & retrieval step.



Persistent Vector DBs

When a large volume of data is stored in vector databases, the retrieval and search needs to be quick. The relevance and accuracy of the search can be tested.



Temporary Vector Index

When data is temporarily stored in vector indices for one time use, the accuracy and relevance of the search needs to be ascertained



Small Data

Generally, when small amount of data is retrieved from pre-determined external sources, the augmentation of the data becomes more critical.

Indexing Pipeline

On the fly

Retrieval

Perhaps, the most critical step in the entire RAG value chain is searching and retrieving the relevant pieces of information (known as **documents**). When the user enters a query or a prompt, it is this system (**Retriever**) that is responsible for accurately fetching the correct snippet of information that is used in responding to the user query.

Retrievers accept a Query as input and return a list of Documents as output

Popular Retrieval Methods

Similarity Search



The similarity search functionality of vector databases forms the backbone of a Retriever. Similarity is calculated by calculating the distance between the embedding vectors of the input and the documents

Maximum Marginal Relevance



MMR addresses redundancy in retrieval. MMR considers the relevance of each document only in terms of how much new information it brings given the previous results. MMR tries to reduce the redundancy of results while at the same time maintaining query relevance of results for already ranked documents/phrases

Multi-query Retrieval



Multi-query Retrieval automates prompt tuning using a language model to generate diverse queries for a user input, retrieving relevant documents from each query and combining them to overcome limitations and obtain a more comprehensive set of results. This approach aims to enhance retrieval performance by considering multiple perspectives on the same query.

Retrieval Methods

Contextual compression

Sometimes, relevant info is hidden in long documents with a lot of extra stuff. Contextual Compression helps with this by squeezing down the documents to only the important parts that match your search.

Multi Vector Retrieval

Sometimes it makes sense to store more than one vectors in a document. E.g A chapter, its summary and a few quotes. The retrieval becomes more efficient because it can match with all the different types of information that has been embedded.

Parent Document Retrieval

In breaking down documents for retrieval, there's a dilemma. Small pieces capture meaning better in embeddings, but if they're too short, context is lost. The Parent Document Retrieval finds a middle ground by storing small chunks. During retrieval, it fetches these bits, then gets the larger documents they came from using their parent IDs

Self Query

A self-querying retriever is a system that can ask itself questions. When you give it a question in normal language, it uses a special process to turn that question into a structured query. Then, it uses this structured query to search through its stored information. This way, it doesn't just compare your question with the documents; it also looks for specific details in the documents based on your question, making the search more efficient and accurate.

Retrieval Methods

Time-weighted Retrieval



This method supplements the semantic similarity search with a time delay. It gives more weightage, then, to documents that are fresher or more used than the ones that are older



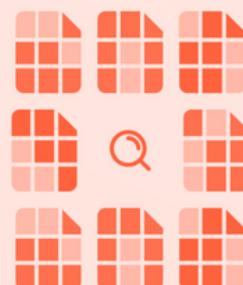
Ensemble Techniques

As the term suggests, multiple retrieval methods can be used in conjunction with each other. There are many ways of implementing ensemble techniques and use cases will define the structure of the retriever

Top Advanced Retrieval Strategies

Top Advanced Retrieval Strategies

- #1 Custom Retrievers
- #2 Self Query
- #3 Hybrid Search
- #4 Contextual Compression
- #5 Multi Query
- #6 TimeWeighted VectorStore



Source : [LangChain State of AI 2023](#)

Example : Similarity Search using LangChain

1. Loading our text file using **TextLoader**,
2. Splitting the text into chunks using **RecursiveCharacterTextSplitter**,
3. Creating embeddings using **all-MiniLM-L6-v2**
4. Storing the embeddings into **Chromadb**
5. Retrieving chunks using **similarity_search**

```
● ● ●  
from langchain.document_loaders import TextLoader  
from langchain.embeddings.sentence_transformer import SentenceTransformerEmbeddings  
from langchain.text_splitter import CharacterTextSplitter  
from langchain.vectorstores import Chroma  
  
# load the document and split it into chunks  
loader = TextLoader('../Data/AK_BusyPersonIntroLLM.txt')  
documents = loader.load()  
  
# split it into chunks  
text_splitter = RecursiveCharacterTextSplitter(chunk_size=1000,  
                                              chunk_overlap=200)  
  
docs = text_splitter.split_documents(documents)  
  
# create the open-source embedding function  
embedding_function = SentenceTransformerEmbeddings(model_name="all-MiniLM-L6-v2")  
  
# load it into Chroma  
db = Chroma.from_documents(docs, embedding_function)  
  
# query it  
query = "What did Andrej say about LLM operating system?"  
docs = db.similarity_search(query)  
  
# print results  
print(docs[0].page_content)
```

```
● ● ●  
...  
llm trying to page relevant information in and out of its  
context window to perform your task um and so a lot of  
other I think connections also exist I think there's  
equivalence of um multi-threading multiprocessing speculative  
execution uh there's equivalent of in the random access memory  
in the context window there's equivalence of user space and  
kernel space and a lot of other equivalents to today's  
operating systems that I didn't fully cover but fundamentally  
the other reason that I really like this analogy of llms kind  
of becoming a bit of an operating system ecosystem is that  
there are also some equivalence I think between the current  
operating systems and the uh and what's emerging today so for  
example in the desktop operating system space we have a few  
proprietary operating systems like Windows and Mac OS but we  
also have this open source ecosystem of a large diversity of  
operating systems based on Linux in the same way here we have  
some proprietary operating systems like GPT  
...  
...
```

Example : Similarity Vector Search

1. Loading our text file using **TextLoader**,
2. Splitting the text into chunks using **RecursiveCharacterTextSplitter**,
3. Creating embeddings using **all-MiniLM-L6-v2**
4. Storing the embeddings into **Chromadb**
5. Converting input query into a **vector embedding**
6. Retrieving chunks using **similarity_search_by_vector**



```
from langchain.document_loaders import TextLoader
from langchain.embeddings.sentence_transformer import SentenceTransformerEmbeddings
from langchain.text_splitter import RecursiveCharacterTextSplitter
from langchain.vectorstores import Chroma

# load the document and split it into chunks
loader = TextLoader('../Data/AK_BusyPersonIntroLLM.txt')
documents = loader.load()

# split it into chunks
text_splitter = RecursiveCharacterTextSplitter(chunk_size=1000,
                                                chunk_overlap=200)

docs = text_splitter.split_documents(documents)

# create the open-source embedding function
embedding_function = SentenceTransformerEmbeddings(model_name="all-MiniLM-L6-v2")

# load it into Chroma
db = Chroma.from_documents(docs, embedding_function)

# query it
query = "What did Andrej say about LLM operating system?"

# convert query to embedding
query_vector=embedding_function.embed_query(query)

# distance based search
docs = db.similarity_search_by_vector(query_vector)

# print results
print(docs[0].page_content)
```



'''llm trying to page relevant information in and out of its context window to perform your task um
and so a lot of other I think connections also exist I think there's equivalence of um multi-threading
multiprocessing speculative execution uh there's equivalent of in the random access memory in the context
window there's equivalence of user space and kernel space and a lot of other equivalents to today's
operating systems that I didn't fully cover but fundamentally the other reason that I really like
this analogy of llms kind of becoming a bit of an operating system ecosystem is that there are also
some equivalence I think between the current operating systems and the uh and what's emerging today
so for example in the desktop operating system space we have a few proprietary operating systems like
Windows and Mac OS but we also have this open source ecosystem of a large diversity of operating systems
based on Linux in the same way here we have some proprietary operating systems like GPT'''

How Similarity Vector Search is different from Similarity Search is that the query is also converted into a vector embedding from regular text

Example : Maximum Marginal Relevance

1. Loading our text file using **TextLoader**,
2. Splitting the text into chunks using **RecursiveCharacterTextSplitter**,
3. Creating embeddings using **OpenAI Embeddings**
4. Storing the embeddings into **Qdrant**
5. Retrieving and ranking chunks using **max_marginal_relevance_search**



```
from langchain.document_loaders import TextLoader
from langchain.embeddings.openai import OpenAIEmbeddings
from langchain.text_splitter import RecursiveCharacterTextSplitter
from langchain.vectorstores import Qdrant

# load the document and split it into chunks
loader = TextLoader('../Data/AK_BusyPersonIntroLLM.txt')
documents = loader.load()

# split it into chunks
text_splitter = RecursiveCharacterTextSplitter(chunk_size=1000,
                                                chunk_overlap=200)

docs = text_splitter.split_documents(documents)

# create the openai embedding function
embedding_function = OpenAIEmbeddings(openai_api_key=openai_api_key)

# load it into Qdrant
db = Qdrant.from_documents(docs, embedding_function, location=":memory:",
                           collection_name="my_documents")

# query it
query = "What did Andrej say about LLM operating system?"

# max marginal relevance search
docs = db.max_marginal_relevance_search(query,k=2, fetch_k=10)

# print results
for i, doc in enumerate(docs):
    print(f"{i + 1}.", doc.page_content, "\n")
```

fetch_k = Number of documents in the initial retrieval
k = final number of reranked documents to output

Example : Multi-query Retrieval

1. Loading our text file using **TextLoader**,
2. Splitting the text into chunks using **RecursiveCharacterTextSplitter**,
3. Creating embeddings using **OpenAI Embeddings**
4. Storing the embeddings into **Qdrant**
5. Set the LLM as **ChatOpenAI (gpt 3.5)**
6. Set up **logging** to see the query variations generated by the LLM
7. use **MultiQueryRetriever & get_relevant_documents** functions



```
from langchain.document_loaders import TextLoader
from langchain.embeddings.openai import OpenAIEMBEDDINGS
from langchain.text_splitter import RecursiveCharacterTextSplitter
from langchain.vectorstores import Qdrant
from langchain.retrievers.multi_query import MultiQueryRetriever
from langchain.chat_models import ChatOpenAI

# load the document and split it into chunks
loader = TextLoader('../Data/AK_BusyPersonIntroLLM.txt')
documents = loader.load()

# split it into chunks
text_splitter = RecursiveCharacterTextSplitter(chunk_size=1000, chunk_overlap=200)
docs = text_splitter.split_documents(documents)

# create the openai embedding function
embedding_function = OpenAIEMBEDDINGS(openai_api_key=openai_api_key)

# load it into Qdrant
db = Qdrant.from_documents(docs, embedding_function, location=":memory:", collection_name="my_documents")

# query it
query = "What did Andrej say about LLM operating system?"

# set the LLM for multiquery
llm = ChatOpenAI(temperature=0, openai_api_key=openai_api_key)

# Multiquery retrieval using OpenAI
retriever_from_llm = MultiQueryRetriever.from_llm(retriever=db.as_retriever(), llm=llm)

# set up logging to see the queries generated
import logging
logging.basicConfig()
logging.getLogger("langchain.retrievers.multi_query").setLevel(logging.INFO)

# retrieved documents
unique_docs = retriever_from_llm.get_relevant_documents(query=query)

# print results
for i, doc in enumerate(unique_docs):
    print(f"{i + 1}.", doc.page_content, "\n")
```

Example : Contextual compression

1. Loading our text file using **TextLoader**,
2. Splitting the text into chunks using **RecursiveCharacterTextSplitter**,
3. Creating embeddings using **OpenAI Embeddings**
4. Set up retriever as **FAISS**
5. Set the LLM as **ChatOpenAI (gpt 3.5)**
6. Use **LLMChainExtractor** as the compressor
7. use **ContextualCompressionRetriever & get_relevant_documents** functions



```
from langchain.document_loaders import TextLoader
from langchain.embeddings.openai import OpenAIEMBEDDINGS
from langchain.text_splitter import RecursiveCharacterTextSplitter
from langchain.vectorstores import FAISS
from langchain.llms import OpenAI
from langchain.retrievers import ContextualCompressionRetriever
from langchain.retrievers.document_compressors import LLMChainExtractor

# load and split text
loader = TextLoader('../Data/AK_BusyPersonIntroLLM.txt')
documents = loader.load()
text_splitter = RecursiveCharacterTextSplitter(chunk_size=1000,
                                               chunk_overlap=200)
docs = text_splitter.split_documents(documents)
# save as vector embeddings
retriever = FAISS.from_documents(docs,
                                  OpenAIEMBEDDINGS(
                                      openai_api_key=openai_api_key)).as_retriever()

# use a compressor
llm = OpenAI(temperature=0, openai_api_key=openai_api_key)
compressor = LLMChainExtractor.from_llm(llm)
compression_retriever = ContextualCompressionRetriever(
    base_compressor=compressor,
    base_retriever=retriever)

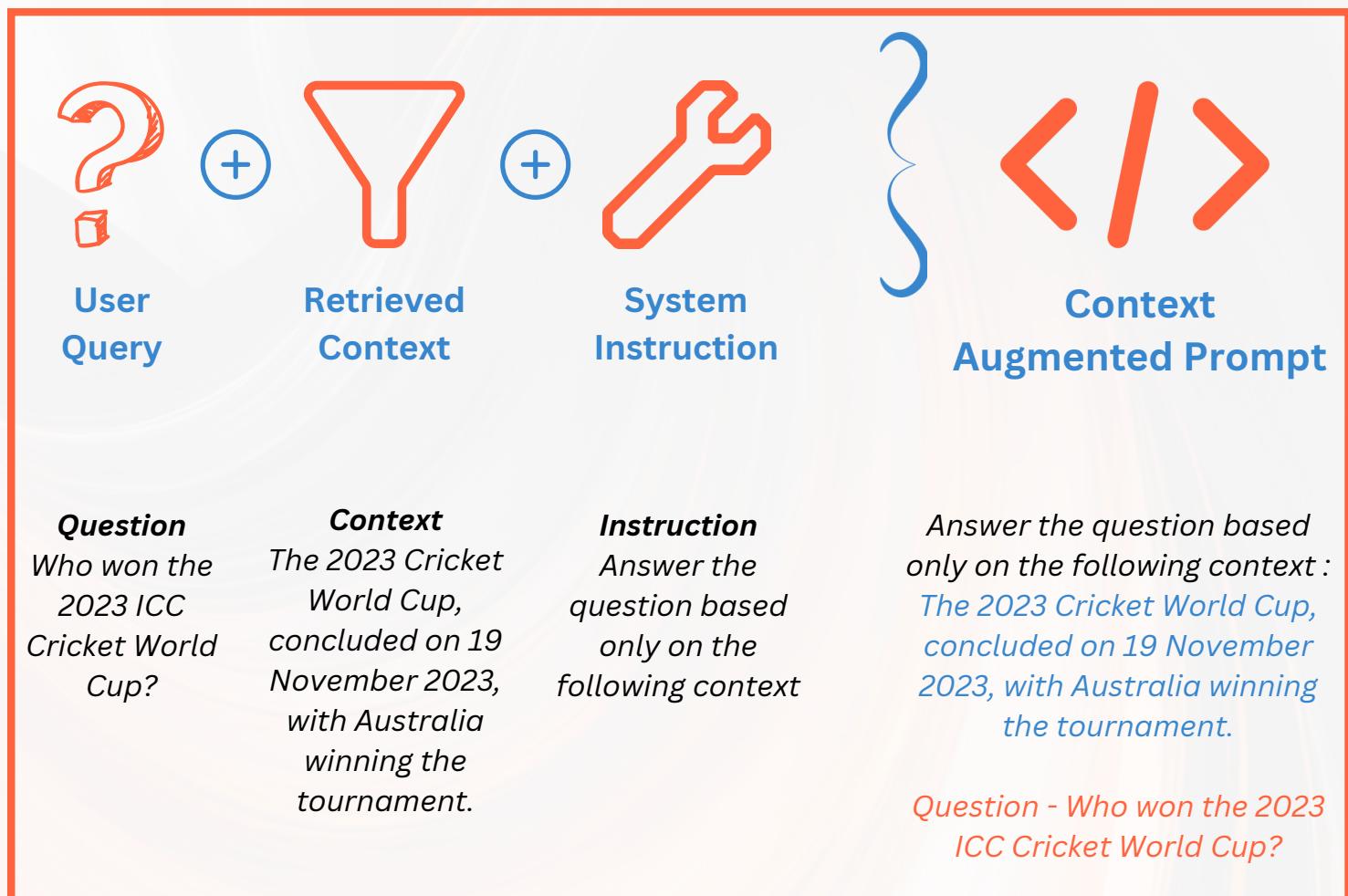
query = "What did Andrej say about LLM operating system?"

# retrieve docs
compressed_docs = compression_retriever.get_relevant_documents(query)

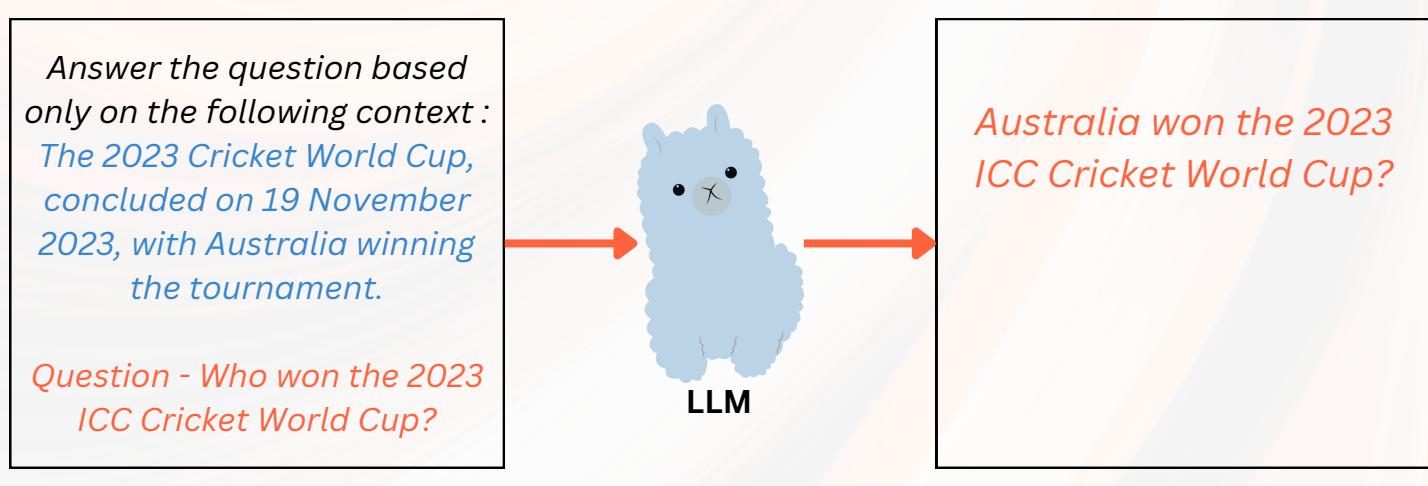
# print docs
for i, doc in enumerate(unique_docs):
    print(f"{i + 1}. {doc.page_content}, \n")
```

Augmentation & Generation

Post-retrieval, the next set of steps include merging the user query and the retrieved context (**Augmentation**) and passing this merged prompt as an instruction to an LLM (**Generation**)



Augmentation with an Illustrative Example



Generation with an Illustrative Example

Evaluation

Building a PoC RAG pipeline is not overtly complex. LangChain and LlamaIndex have made it quite simple. Developing highly impressive Large Language Model (LLM) applications is achievable through brief training and verification on a limited set of examples. However, to enhance its robustness, thorough testing on a dataset that accurately mirrors the production distribution is imperative.

RAG is a great tool to address hallucinations in LLMs but...
even RAGs can suffer from hallucinations

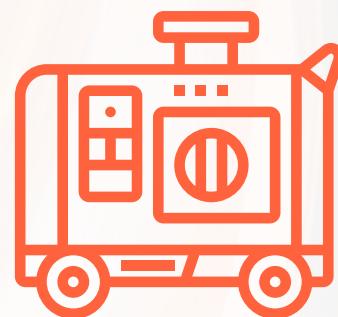
This can be because -

- The retriever fails to retrieve relevant context or retrieves irrelevant context
- The LLM, despite being provided the context, does not consider it
- The LLM instead of answering the query picks irrelevant information from the context

Two processes, therefore, to focus on from an evaluation perspective -



Search & Retrieval



Generation



How good is the retrieval of the context from the Vector Database?



Is it relevant to the query?



How much noise (irrelevant information) is present?



How good is the generated response?



Is the response grounded in the provided context?



Is the response relevant to the query?

Ragas (RAG Assessment)

Jithin James and Shahul ES from Exploding Gradients, in 2023, developed the Ragas framework to address these questions.

<https://github.com/explodinggradients/ragas> 

Evaluation Data

To evaluate RAG pipelines, the following four data points are recommended

-  A set of **Queries** or **Prompts** for evaluation
-  **Retrieved Context** for each prompt
-  Corresponding **Response** or **Answer** from LLM
-  **Ground Truth** or known correct response

Evaluation Metrics

Evaluating Generation

-  **Faithfulness** Is the **Response** faithful to the **Retrieved Context**?
- Answer Relevance** Is the **Response** relevant to the **Prompt**?

Retrieval Evaluation

-  **Context Relevance** Is the **Retrieved Context** relevant to the **Prompt**?
- Context Recall** Is the **Retrieved Context** aligned to the **Ground Truth**?
- Context Precision** Is the **Retrieved Context** ordered correctly?

Overall Evaluation

- Answer Semantic Similarity** Is the **Response** semantically similar to the **Ground Truth**?

- Answer Correctness** Is the **Response** semantically and factually similar to the **Ground Truth**?

Evaluation Metrics

1 Faithfulness

Faithfulness is the measure of the extent to which the response is factually grounded in the retrieved context

Problem addressed : The LLM, despite being provided the context, does not consider it

or

Is the response grounded in the provided context?

Evaluated Process : Generation

Any measure of retrieval accuracy is out of scope

Score Range : (0,1) Higher score is better

Methodology

Faithfulness identifies the number of “claims” made in the response and calculates the proportion of those “claims” present in the context.

$$\text{Faithfulness} = \frac{\text{Number of generated claims present in the context}}{\text{Total number of claims made in the generated response}}$$

Illustrative Example

Query : Who won the 2023 ODI Cricket World Cup and when?

Context : The 2023 ODI Cricket World Cup concluded on 19 November 2023, with Australia winning the tournament.

Response 1 : High Faithfulness

[Australia] won on [19 November 2023]

Response 2 : Low Faithfulness

[Australia] won on [15 October 2023]

Evaluation Metrics

2 Answer Relevance

Answer Relevance is the measure of the extent to which the response is relevant to the query or the prompt

Problem addressed : The LLM instead of answering the query responds with irrelevant information

or

Is the response relevant to the query?

Evaluated Process : Generation

Any measure of retrieval accuracy is out of scope

Score Range : (0,1) **Higher score is better**

Methodology

For this metric, a response is generated for the initial query or prompt. To compute the score, the LLM is then prompted to generate questions for the generated response several times. The mean cosine similarity between these questions and the original one is then calculated. The concept is that if the answer correctly addresses the initial question, the LLM should generate questions from it that match the original question.

Avg (

Answer Relevance = *Sc (Initial Query, LLM generated Query [i])*
)

Illustrative Example

Query : Who won the 2023 ODI Cricket World Cup and when?

Response 1 : High Answer Relevance

India won on 19 November 2023

Response 2 : Low Answer Relevance

Cricket world cup is held once every four years

Note

Answer Relevance is **not a measure of truthfulness** but only of relevance. The response may or may not be factually accurate but may be relevant.

Evaluation Metrics

3 Context Relevance

Context Relevance is the measure of the extent to which the retrieved context is relevant to the query or the prompt

Problem addressed : The retriever fails to retrieve relevant context
or

Is the retrieved context relevant to the query?

Evaluated Process : Retrieval

Indifferent to the final generated response

Score Range : (0,1) Higher score is better

Methodology

The retrieved context should contain information only relevant to the query or the prompt. For context relevance, a metric 'S' is estimated. 'S' is the number of sentences in the retrieved context that are relevant for responding to the query or the prompt.

$$\text{Context Relevance} = \frac{S}{\text{(number of relevant sentences from the context)}} \\ \text{Total number of sentences in the retrieved context}$$

Illustrative Example

Query : Who won the 2023 ODI Cricket World Cup and when?

Context 1 : High Context Relevance

The 2023 Cricket World Cup, concluded on 19 November 2023, with Australia winning the tournament. The tournament took place in ten different stadiums, in ten cities across the country. The final took place between India and Australia at Narendra Modi Stadium

Context 2 : Low Context Relevance

The 2023 Cricket World Cup was the 13th edition of the Cricket World Cup. It was the first Cricket World Cup which India hosted solely. The tournament took place in ten different stadiums. In the first semi-final India beat New Zealand, and in the second semi-final Australia beat South Africa.

Evaluation Metrics

Ground Truth

Ground truth is information that is known to be real or true. In RAG, or Generative AI domain in general, Ground Truth is a prepared set of **Prompt-Response examples**. It is akin to *labelled data* in Supervised Learning parlance.

Calculation of certain metrics necessitates the availability of Ground Truth data

4

Context Recall

Context recall measures the extent to which the retrieved context aligns with the “provided” answer or Ground Truth

Problem addressed : The retriever fails to retrieve accurate context
or

Is the retrieved context good enough to provide the response?

Evaluated Process : Retrieval

Indifferent to the final generated response

Score Range : (0,1) Higher score is better

Methodology

To estimate context recall from the ground truth answer, each sentence in the ground truth answer is analyzed to determine whether it can be attributed to the retrieved context or not. Ideally, all sentences in the ground truth answer should be attributable to the retrieved context.

Context Recall =

$$\frac{\text{Number of Ground Truth sentences in the context}}{\text{Total number of sentences in the Ground Truth}}$$

Illustrative Example

Query : Who won the 2023 ODI Cricket World Cup and when?

Ground Truth : Australia won the world cup on 19 November, 2023.

Context 1 : High Context Recall

The 2023 Cricket World Cup, concluded on 19 November 2023, with Australia winning the tournament.

Context 2 : Low Context Recall

The 2023 Cricket World Cup was the 13th edition of the Cricket World Cup. It was the first Cricket World Cup which India hosted solely.

Evaluation Metrics

5 Context Precision

Context Precision is a metric that evaluates whether all of the ground-truth relevant items present in the contexts are ranked higher or not.

Problem addressed : The retriever fails to rank retrieve context correctly
or

Is the higher ranked retrieved context better to provide the response?

Evaluated Process : Retrieval

Indifferent to the final generated response

Score Range : (0,1) **Higher score is better**

Methodology

Context Precision is a metric that evaluates whether all of the ground-truth relevant items present in the all retrieved context documents are ranked higher or not. Ideally all the relevant chunks must appear at the top

$$\text{Context Precision @ k} = \frac{\text{Sum(Precision@k)}}{\text{Total number of relevant documents in the top results}}$$

$$\text{Precision @ k} = \frac{\text{True Positives @ k}}{(\text{True Positives @ k} + \text{False Positives @ k})}$$

Precision @ k

Precision@k is a metric used in information retrieval and recommendation systems to evaluate the accuracy of the top k items retrieved or recommended. It measures the proportion of relevant items among the top k items.

Evaluation Metrics

6 Answer semantic similarity

Answer semantic similarity evaluates whether the generated response is similar to the “provided” response or Ground Truth.

Problem addressed : The generated response is incorrect
or

Does the pipeline generate the right response?

Evaluated Process : Retrieval & Generation

Score Range : (0,1) Higher score is better

Methodology

Answer semantic similarity score is calculated by measuring the semantic similarity between the generated response and the ground truth response.

$$\text{Answer Semantic Similarity} = \frac{\text{Similarity}(\text{Generated Response}, \text{Ground Truth Response})}{\text{Similarity}}$$

7 Answer Correctness

Answer correctness evaluates whether the generated response is semantically and factually similar to the “provided” response or Ground Truth.

Problem addressed : The generated response is incorrect
or

Does the pipeline generate the right response?

Evaluated Process : Retrieval & Generation

Score Range : (0,1) Higher score is better

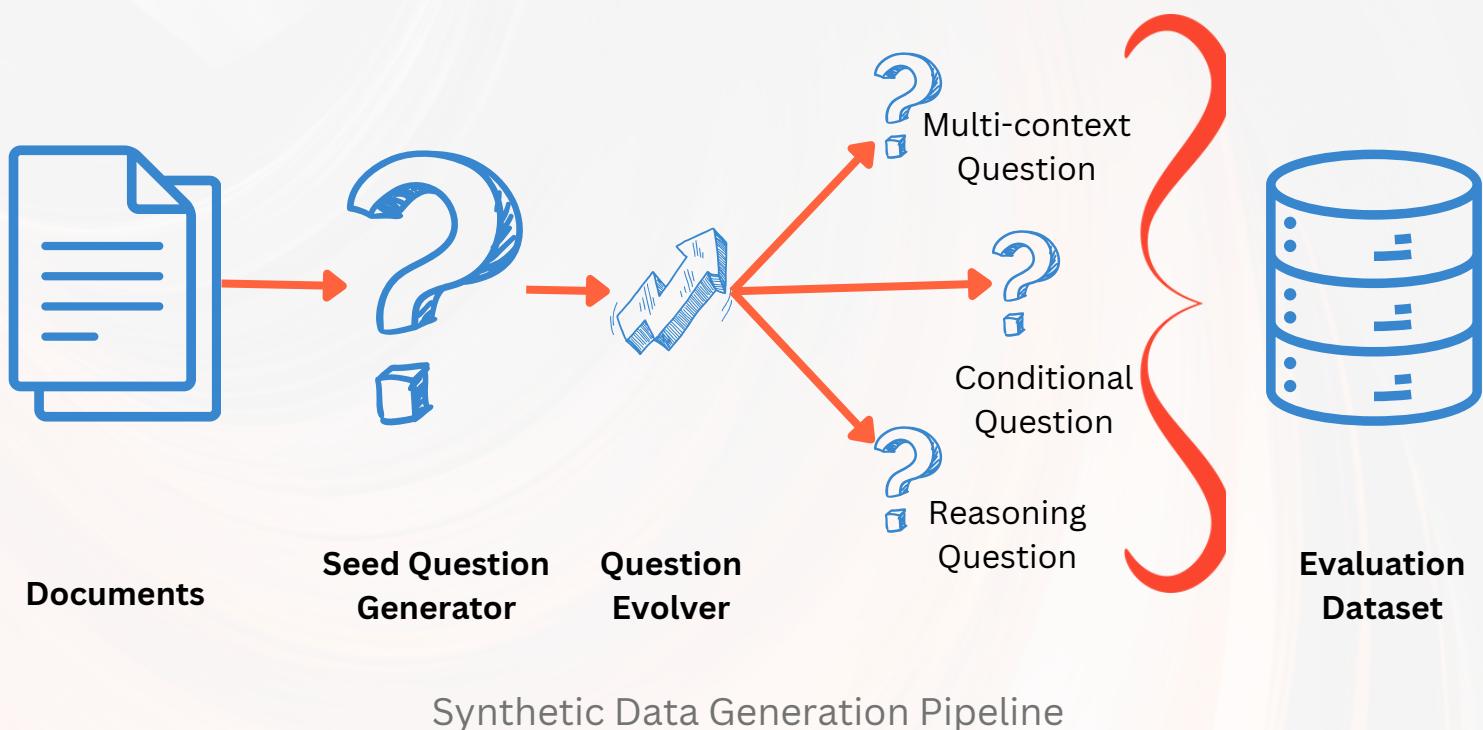
Methodology

Answer correctness score is calculated by measuring the semantic and the factual similarity between the generated response and the ground truth response.

Synthetic Test Data Generation

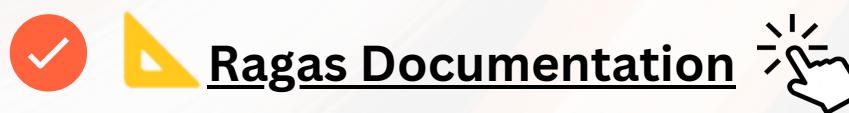
Generating hundreds of QA (Question-Context-Answer) samples from documents manually can be a time-consuming and labor-intensive task. Moreover, questions created by humans may face challenges in achieving the necessary level of complexity for a comprehensive evaluation, potentially affecting the overall quality of the assessment.

Synthetic Data Generation uses Large Language Models to generate a variety of Questions/Prompts and Responses/Answers from the Documents (Context). It can greatly reduce developer time.



	question	context	answer	question_type	episode_done
0	What technique improves the performance of lar...	- "We explore how generating a chain of thought..."	The technique that improves the performance of...	simple	True
1	What phenomenon is discussed in the paper rega...	- This paper instead discusses an unpredictabl...	The phenomenon discussed in the paper is the e...	reasoning	True
2	What is the purpose of chain-of-thought (CoT) ...	- Providing these steps for prompting demonstr...	The purpose of chain-of-thought (CoT) promptin...	simple	True
3	What is the performance of the largest fine-tu...	On the MathQA-Python dataset, the largest fine...	The performance of the largest fine-tuned mode...	simple	True
4	What is the accuracy increase of Zero-shot-CoT...	Experimental results demonstrate that our Zero...	The accuracy increase of Zero-shot-CoT on Mult...	reasoning	True

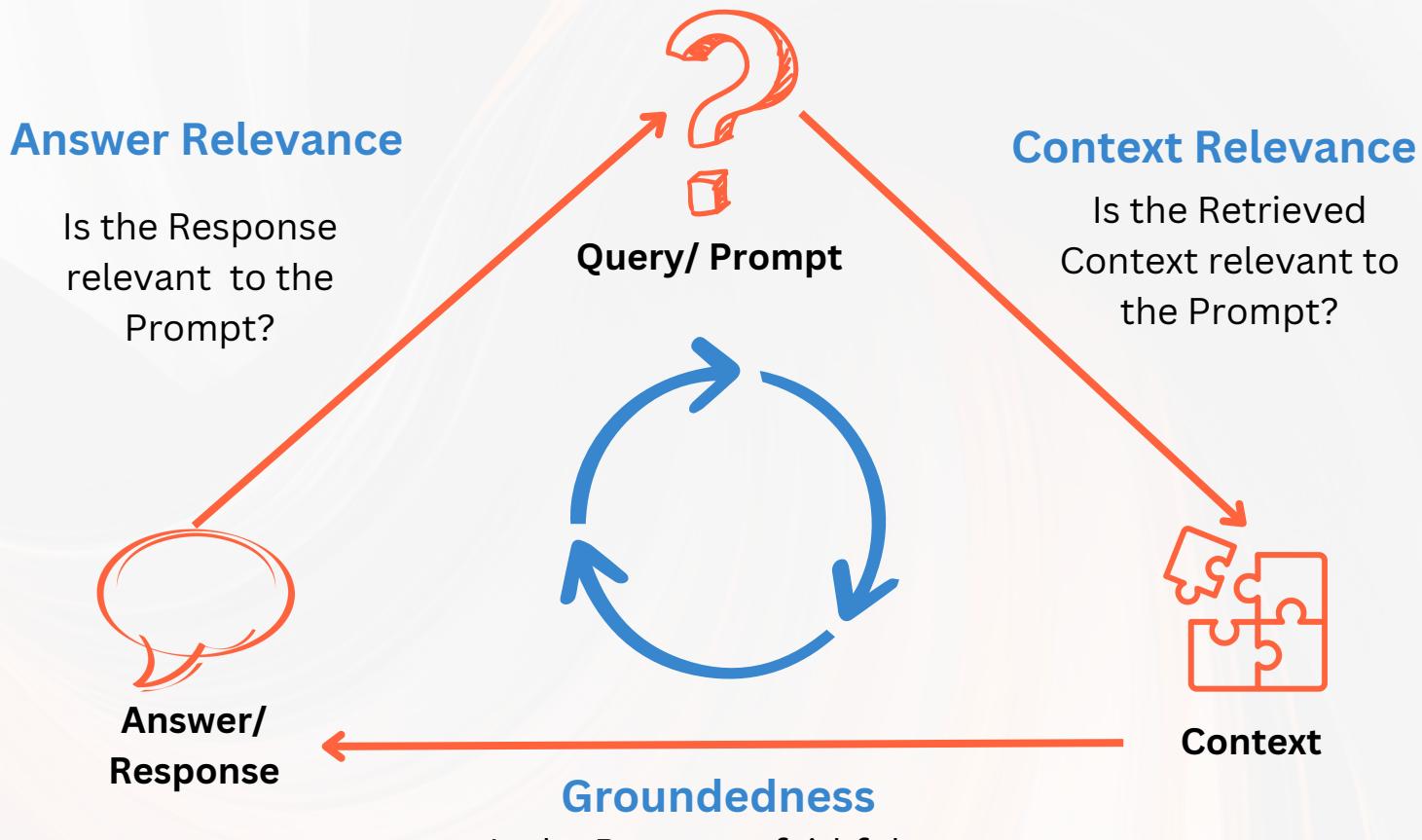
Synthetic Data Generated Using Ragas





The RAG Triad (TruLens)

The RAG triad is a framework proposed by TruLens to evaluate hallucinations along each edge of the RAG architecture.



Context Relevance:

- Verify quality by ensuring each context chunk is relevant to the input query

Groundedness:

- Verify groundedness by breaking down the response into individual claims.
- Independently search for evidence supporting each claim in the retrieved context.

Answer Relevance:

- Ensure the response effectively addresses the original question.
- Verify by evaluating the relevance of the final response to user input.



[Trulens Documentation](#)



RAG vs Finetuning vs Both

Supervised Finetuning (SFT) has fast become a popular method to customise and adapt foundation models for specific objectives. There has been a growing debate in the applied AI community around the application of fine-tuning or RAG to accomplish tasks.

RAG & SFT should be considered as complementary, rather than competing, techniques.

RAG enhances the non-parametric memory of a foundation model without changing the parameters

SFT changes the parameters of a foundation model and therefore impacting the parametric memory

If the requirement dictates changes to the parametric memory and an increase in the non-parametric memory, then RAG and SFT can be used in conjunction

RAG Features

Connect to dynamic external data sources ✓

Reduce hallucinations ✓

Increase transparency (in terms of source of information) ✓

Works well only with very large foundation models ✗

Does not impact the style, tone, vocabulary of the foundation model ✗

SFT Features

Change the style, vocabulary, tone of the foundation model ✓

Can reduce model size ✓

Useful for deep domain expertise ✓

May not address the problem of hallucinations ✗

No improvement in transparency (as black box as foundation models) ✗

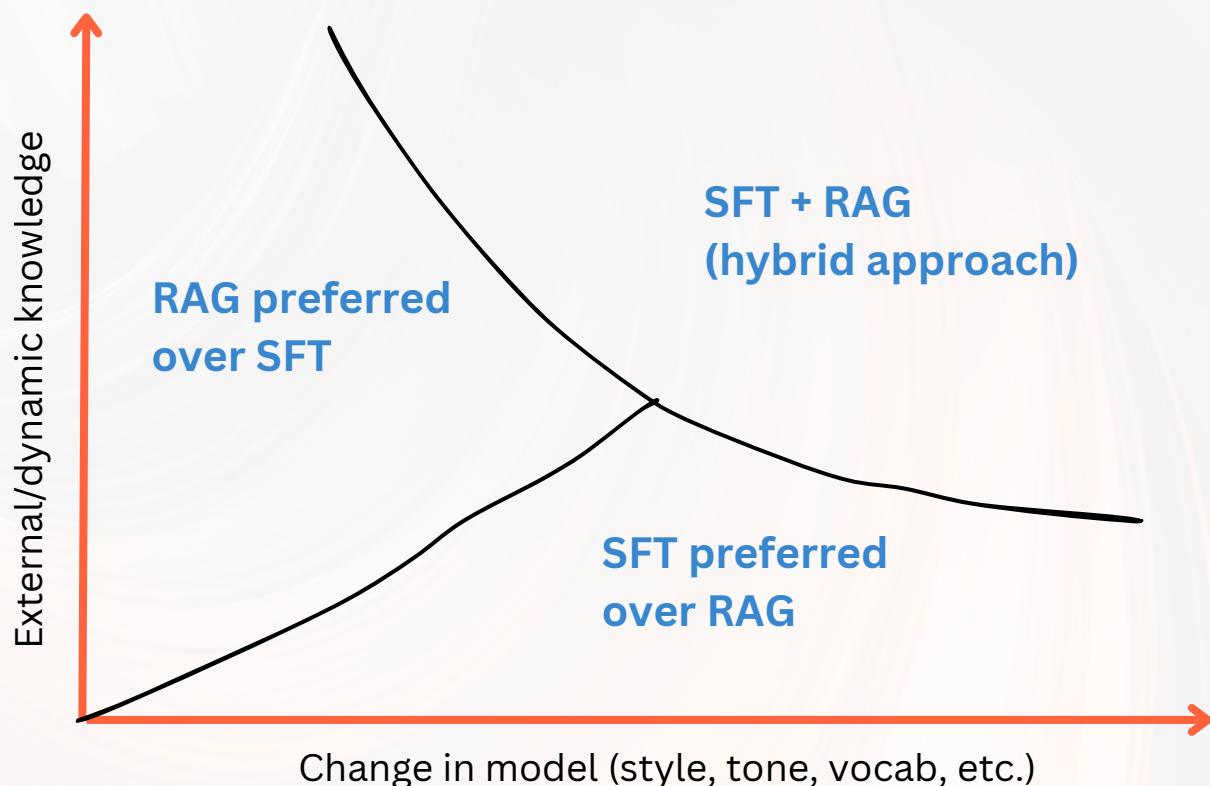
Important Use Case Considerations

Do you require usage of dynamic external data?

RAG preferred over SFT

Do you require changing the writing style, tonality, vocabulary of the model?

SFT preferred over RAG



RAG should be implemented (with or without SFT) if the use case requires

- Access to an external data source, especially, if the data is dynamic
- Resolving Hallucinations
- Transparency in terms of the source of information

For SFT, you'll need to have access to labelled training data

Other Considerations

Latency

RAG pipelines require an additional step of searching and retrieving context which introduces an inherent latency in the system

Scalability

RAG pipelines are modular and therefore can be scaled relatively easily when compared to SFT. SFT will require retraining the model with each additional data source

Cost

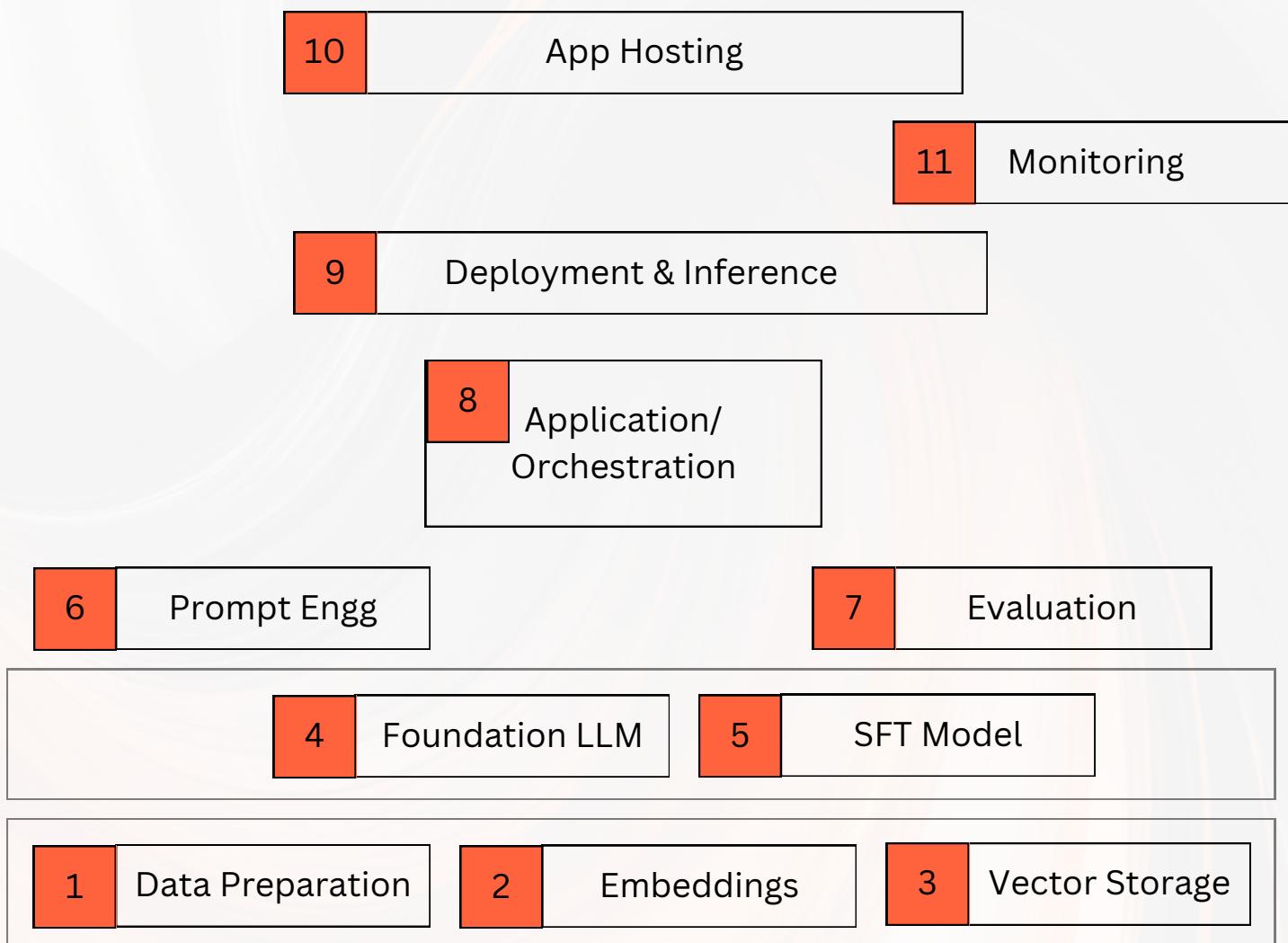
Both RAG and SFT warrant upfront investment. Training cost for SFT can vary depending on the technique and the choice of foundation model. Setting up the knowledge base and integration can be costly for RAG

Expertise

Creating RAG pipelines has become moderately simple with frameworks like LangChain and LlamaIndex. Fine-tuning on the other hand requires deep understanding of the techniques and creation of training data

Evolving RAG LLMops Stack

The production ecosystem for RAG and LLM applications is still evolving. Early tooling and design patterns have emerged.

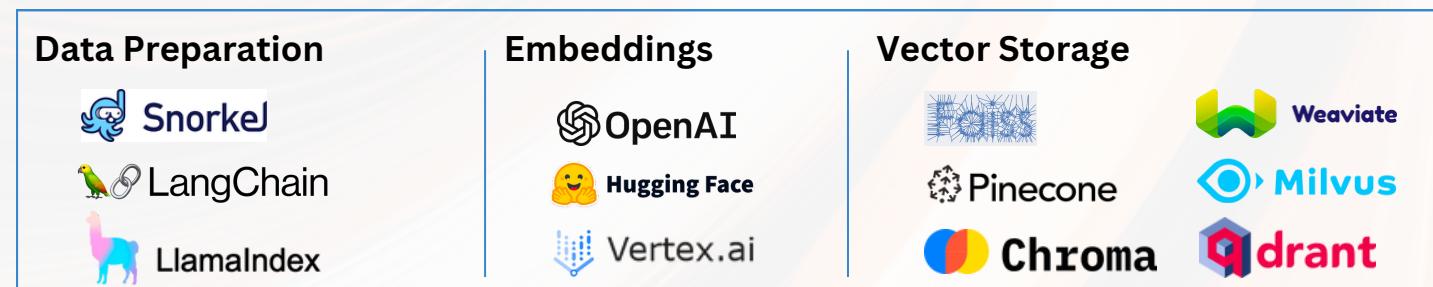


Data Layer

The foundation of RAG applications is the data layer. This involves -

- Data preparation - Sourcing, Cleaning, Loading & Chunking
- Creation of Embeddings
- Storing the embeddings in a vector store

We've seen this process in the creation of the **indexing pipeline**



Popular Data Layer Vendors (Non Exhaustive)

Model Layer

2023 can be considered a year of LLM wars. Almost every other week in the second half of the year a new model was released. Like there is no RAG without data, there is no RAG without an LLM. There are four broad categories of LLMs that can be a part of a RAG application

- 1. A Proprietary Foundation Model** - Developed and maintained by providers (like OpenAI, Anthropic, Google) and is generally available via an API
- 2. Open Source Foundation Model** - Available in public domain (like Falcon, Llama, Mistral) and has to be hosted and maintained by you.
- 3. A Supervised Fine-Tuned Proprietary Model** - Providers enable fine-tuning of their proprietary models with your data. The fine-tuned models are still hosted and maintained by the providers and are available via an API
- 4. A Supervised Fine-Tuned Open Source Model** - All Open Source models can be fine-tuned by you on your data using full fine-tuning or PEFT methods.

There are a lot of vendors that have enabled access to open source models and also facilitate easy fine tuning of these models

Proprietary Models



ANTHROPIC Claude



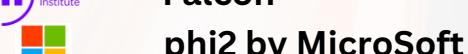
Open Source Models



Mistral & Mixtral

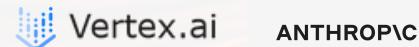


Falcon



Popular proprietary and open source LLMs (Non Exhaustive)

Proprietary Models



ANTHROPIC

Open Source Models



Popular vendors providing access to LLMs (Non Exhaustive)

Note : For Open Source models it is important to check the license type. Some open source models are not available for commercial use

Prompt Layer

Prompt Engineering is more than writing questions in natural language. There are several prompting techniques and developers need to create prompts tailored to the use cases. This process often involves experimentation: the developer creates a prompt, observes the results and then iterates on the prompts to improve the effectiveness of the app. This requires tracking and collaboration



Popular prompt engineering platforms (Non Exhaustive)

Evaluation

It is easy to build a RAG pipeline but to get it ready for production involves robust evaluation of the performance of the pipeline. For checking hallucinations, relevance and accuracy there are several frameworks and tools that have come up.



Popular RAG evaluation frameworks and tools (Non Exhaustive)

App Orchestration

An RAG application involves interaction of multiple tools and services. To run the RAG pipeline, a solid orchestration framework is required that invokes these different processes.

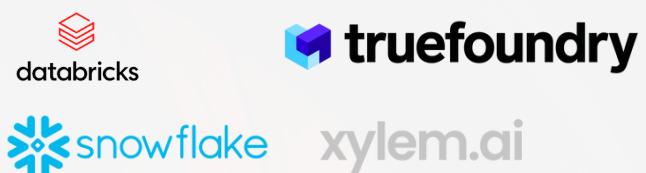


Popular App orchestration frameworks (Non Exhaustive)

Deployment Layer

Deployment of the RAG application can be done on any of the available cloud providers and platforms. Some important factors to consider while deployment are also -

- Security and Governance
- Logging
- Inference costs and latency



Popular cloud providers and LLMops platforms (Non Exhaustive)

Application Layer

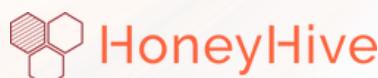
The application finally needs to be hosted for the intended users or systems to interact with it. You can create your own application layer or use the available platforms.



Popular app hosting platforms (Non Exhaustive)

Monitoring

Deployed application needs to be continuously monitored for both accuracy and relevance as well as cost and latency.



Popular monitoring platforms (Non Exhaustive)

Other Considerations

LLM Cache - To reduce costs by saving responses for popular queries

LLM Guardrails - To add additional layer of scrutiny on generations

Multimodal RAG

Up until now, most AI models have been limited to a single modality (a single type of data like text or images or video). Recently, there has been significant progress in AI models being able to handle multiple modalities (majorly text and images). With the emergence of these Large Multimodal Models (LMMs) a multimodal RAG system becomes possible.

“Generate any type of output from any type of input providing any type of context”

The high-level features of multimodal RAG are -

1. Ability to **query/prompt in one or more modalities** like sending both text and image as input.
2. Ability to **search and retrieve not only text** but also images, tables, audio files related to the query
3. Ability to **generate text, image, video etc.** irrespective of the mode(s) in which the input is provided.

Approaches



Using MultiModal Embeddings



Using LMMs Only

Large MultiModel Models



Flamingo



BLIP



KOSMOS-1



Macaw-LLM



GPT4



Gemini

LlaVA

LAVIN

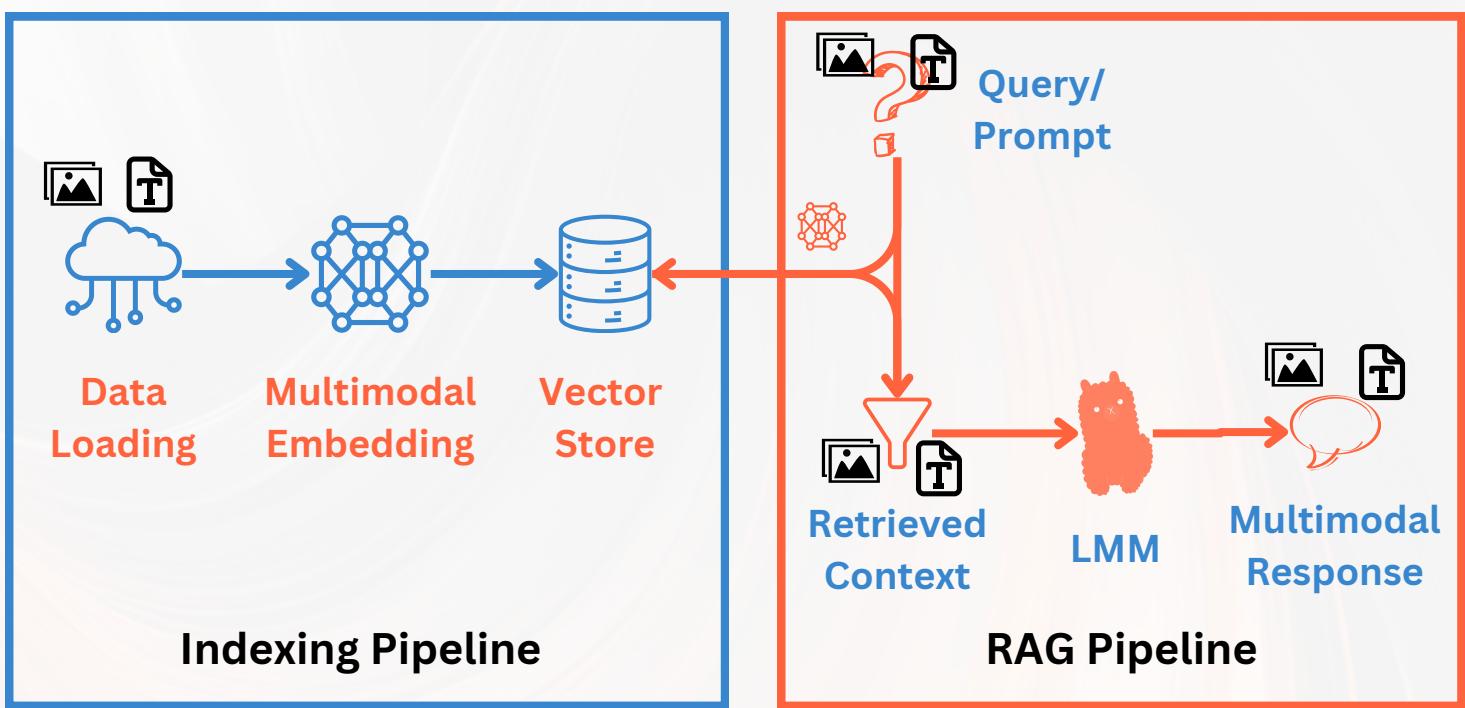
LLaMA - Adapter

FUYU

Multimodal RAG Approaches

Using MultiModal Embeddings

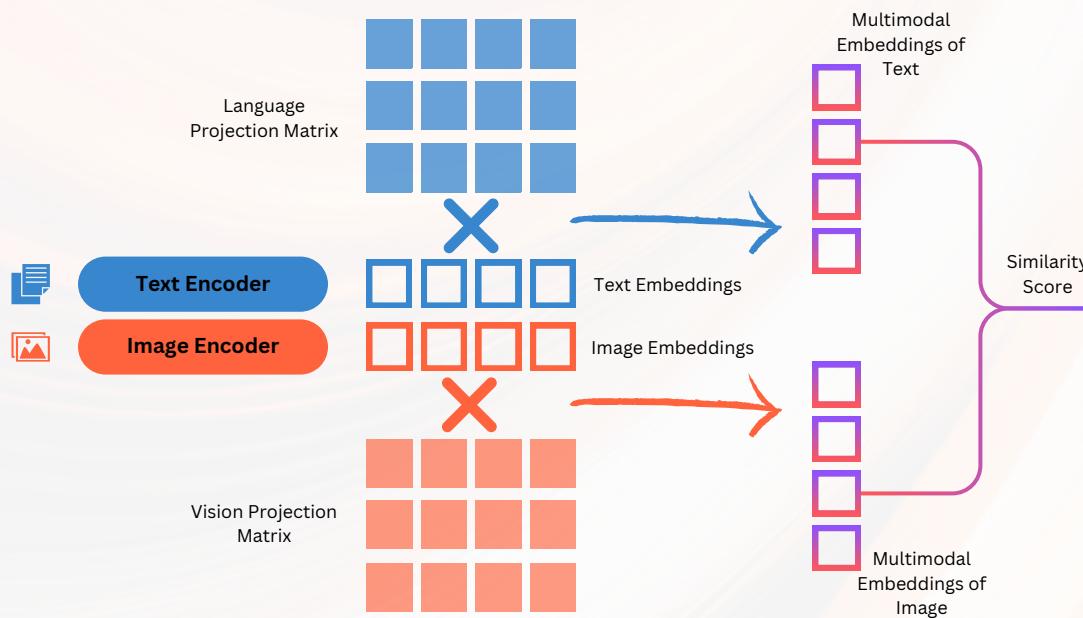
- Multimodal embeddings (like **CLIP**) are used to embed images and text
- User Query is used to retrieve context which can be image and/or text
- The image and/or text context is passed to an LMM with the prompt.
- The LMM generates the final response based on the prompt



Multimodal RAG using Multimodal Embeddings

CLIP : Contrastive Language-Image Pre-training

Mapping data of different modalities into a shared embedding space



CLIP is an example of training multimodal embeddings

OpenAI's **CLIP (Contrastive Language-Image Pre-training)**, maps both images and text into the same semantic embedding space. This allows CLIP to "understand" the relationship between texts and images for powerful applications

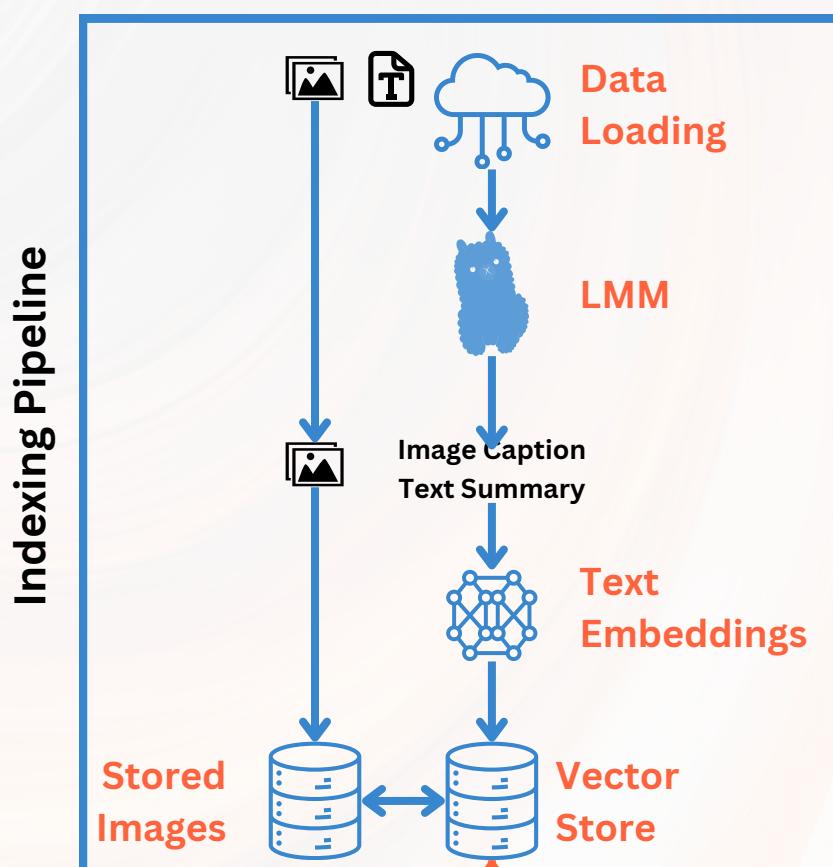
Using LMMs to produce text summaries from images

Indexing

- An LLM is used to generate captions for images in the data
- The image captions and text summaries are stored as text embeddings in a vector database
- A mapping is maintained from the image captions to the image files

Generation

- User enters a query (with text and image)
- Image captions are generated using an LLM and embeddings are generated
- Text summaries and image captions are searched. Images are retrieved based on the relevant image captions.
- Retrieved text summaries, captions and images are passed to the LMM with the prompt. The LMM generates a multimodal response



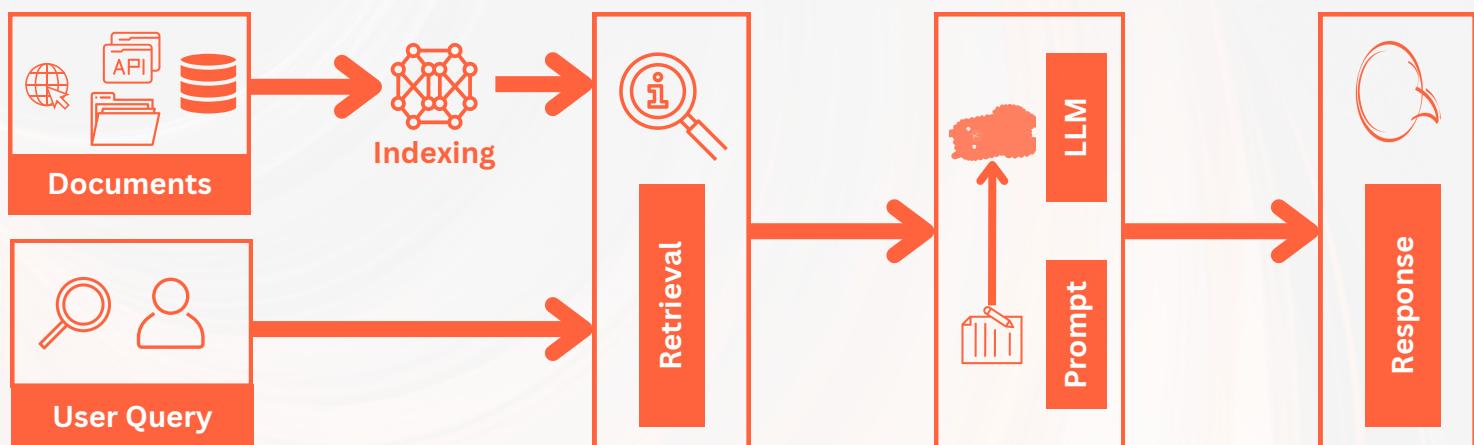
Progression of RAG Systems

Ever since its introduction in mid-2020, RAG approaches have followed a progression aiming to achieve the redressal of the hallucination problem in LLMs

Naive RAG

At its most basic, Retrieval Augmented Generation can be summarized in three steps -

1. **Indexing** of the **documents**
2. **Retrieval** of the context with respect to an input query
3. **Generation** of the **response** using the input query and retrieved context



This basic RAG approach can also be termed “Naive RAG”

Challenges in Naive RAG

Retrieval Quality

- **Low Precision** leading to Hallucinations/Mid-air drops
- **Low Recall** resulting in missing relevant info
- **Outdated information**

Augmentation

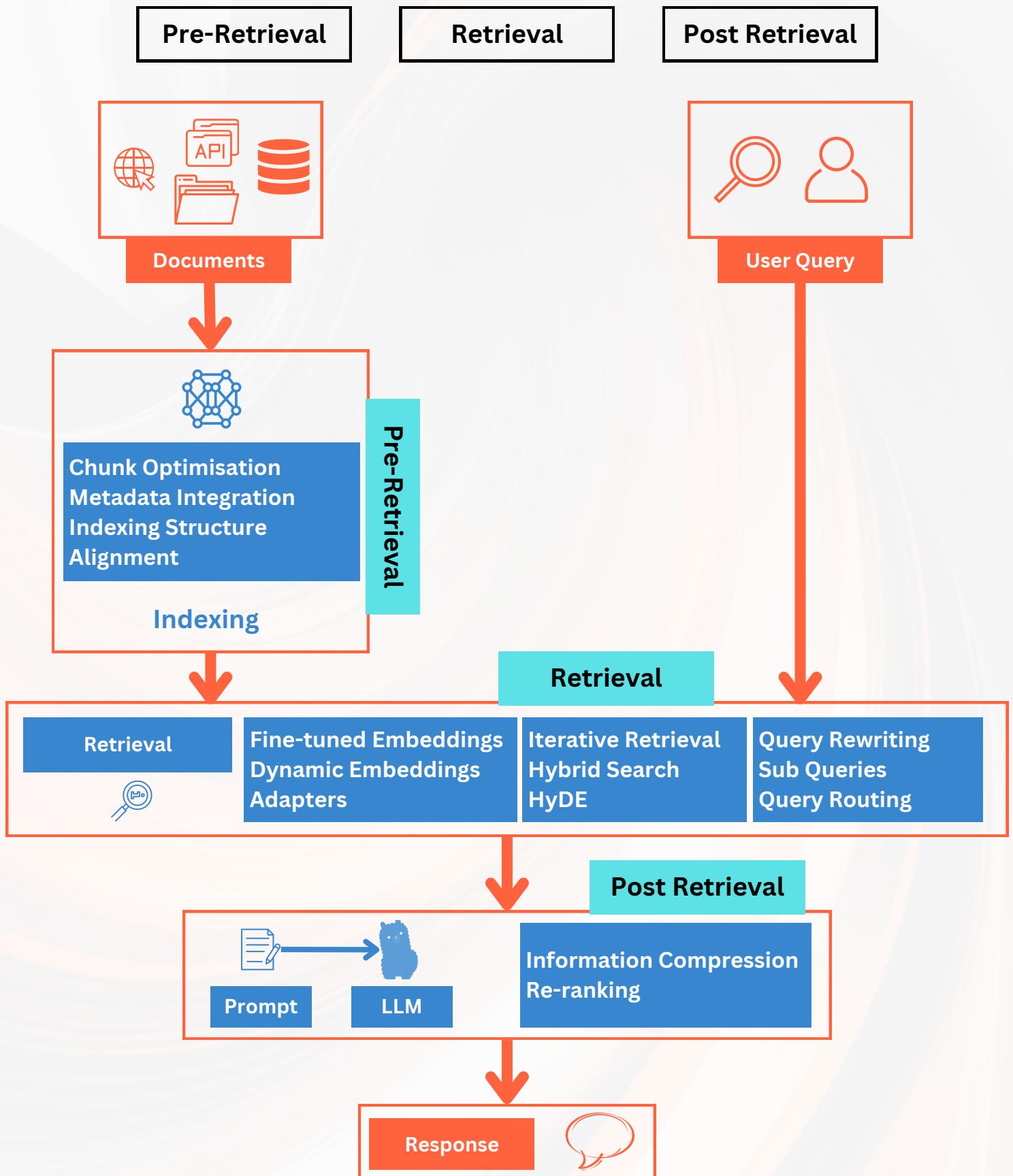
- **Redundancy and Repetition** when multiple retrieved documents have similar information
- **Context Length** challenges

Generation Quality

- Generations are **not grounded** in the context
- Potential of **toxicity and bias** in the response
- **Excessive dependence** on augmented context

Advanced RAG

To address the inefficiencies of the Naive RAG approach, Advanced RAG approaches implement strategies focussed on three processes -



* Indicative, non-exhaustive list

Advanced RAG Concepts

Pre-retrieval/Retrieval Stage

Chunk Optimization

When managing external documents, it's important to break them into the right-sized chunks for accurate results. The choice of how to do this depends on factors like content type, user queries, and application needs. No one-size-fits-all strategy exists, so flexibility is crucial. Current research explores techniques like sliding windows and "small2big" methods.

Metadata Integration

Information like dates, purpose, chapter summaries, etc. can be embedded into chunks. This improves the retriever efficiency by not only searching the documents but also by assessing the similarity to the metadata.

Indexing Structure

Introduction of graph structures can greatly enhance retrieval by leveraging nodes and their relationships. Multi-index paths can be created aimed at increasing efficiency.

Alignment

Understanding complex data, like tables, can be tricky for RAG. One way to improve the indexing is by using counterfactual training, where we create hypothetical (what-if) questions. This increases the alignment and reduces disparity between documents.

Query Rewriting

To bring better alignment between the user query and documents, several rewriting approaches exist. LLMs are sometimes used to create pseudo documents from the query for better matching with existing documents. Sometimes, LLMs perform abstract reasoning. Multi-querying is employed to solve complex user queries.

Hybrid Search Exploration

The RAG system employs different types of searches like keyword, semantic and vector search, depending upon the user query and the type of data available.

Sub Queries

Sub querying involves breaking down a complex query into sub questions for each relevant data source, then gather all the intermediate responses and synthesize a final response.

Query Routing

A query router identifies a downstream task and decides the subsequent action that the RAG system should take. During retrieval, the query router also identifies the most appropriate data source for resolving the query.

Iterative Retrieval

Documents are collected repeatedly based on the query and the generated response to create a more comprehensive knowledge base.

Recursive Retrieval

Recursive retrieval also iteratively retrieves documents. However, it also refines the search queries depending on the results obtained from the previous retrieval. It is like a continuous learning process.

Adaptive Retrieval

Enhance the RAG framework by empowering Language Models (LLMs) to proactively identify the most suitable moments and content for retrieval. This refinement aims to improve the efficiency and relevance of the information obtained, allowing the models to dynamically choose when and what to retrieve, leading to more precise and effective results

Hypothetical Document Embeddings (HyDE)

Using the Language Model (LLM), HyDE forms a hypothetical document (answer) in response to a query, embeds it, and then retrieves real documents similar to this hypothetical one. Instead of relying on embedding similarity based on the query, it emphasizes the similarity between embeddings of different answers.

Fine-tuned Embeddings

This process involves tailoring embedding models to improve retrieval accuracy, particularly in specialized domains dealing with uncommon or evolving terms. The fine-tuning process utilizes training data generated with language models where questions grounded in document chunks are generated.

Post Retrieval Stage

Information Compression

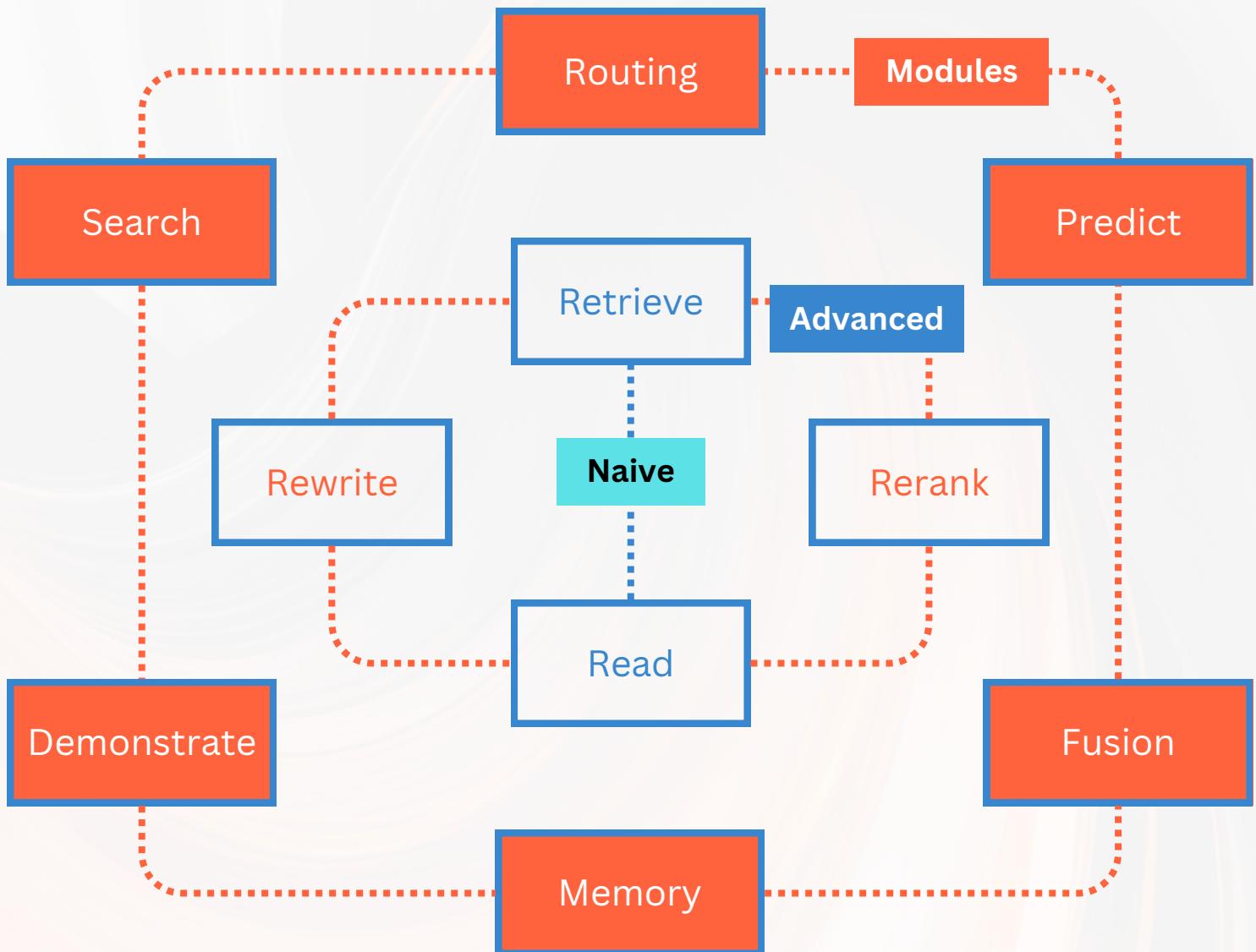
While the retriever is proficient in extracting relevant information from extensive knowledge bases, managing the vast amount of information within retrieval documents poses a challenge. The retrieved information is compressed to extract the most relevant points before passing it to the LLM.

Reranking

The re-ranking model plays a crucial role in optimizing the document set retrieved by the retriever. The main idea is to rearrange document records to prioritize the most relevant ones at the top, effectively managing the total number of documents. This not only resolves challenges related to context window expansion during retrieval but also improves efficiency and responsiveness.

Modular RAG

The SOTA in Retrieval Augmented Generation is a modular approach which allows components like search, memory, and reranking modules to be configured



Naive RAG is essentially a **Retrieve -> Read** approach which focusses on retrieving information and comprehending it.

Advanced RAG adds to the **Retrieve -> Read** approach by adding it into a **Rewrite** and **Rerank** components to improve relevance and groundedness.

Modular RAG takes everything a notch ahead by providing **flexibility** and adding modules like **Search, Routing, etc.**

Naive, Advanced & Modular RAGs are **not exclusive approaches** but a **progression**. Naive RAG is a special case of Advanced which, in turn, is a special case of Modular RAG

Some RAG Modules

Search

The search module is aimed at performing search on different data sources. It is customised to different data sources and aimed at increasing the source data for better response generation

Memory

This module leverages the parametric memory capabilities of the Language Model (LLM) to guide retrieval. The module may use a retrieval-enhanced generator to create an unbounded memory pool iteratively, combining the "original question" and "dual question." By employing a retrieval-enhanced generative model that improves itself using its own outputs, the text becomes more aligned with the data distribution during the reasoning process.

Fusion

RAG-Fusion improves traditional search systems by overcoming their limitations through a multi-query approach. It expands user queries into multiple diverse perspectives using a Language Model (LLM). This strategy goes beyond capturing explicit information and delves into uncovering deeper, transformative knowledge. The fusion process involves conducting parallel vector searches for both the original and expanded queries, intelligently re-ranking to optimize results, and pairing the best outcomes with new queries.

Extra Generation

Rather than directly fetching information from a data source, this module employs the Language Model (LLM) to generate the required context. The content produced by the LLM is more likely to contain pertinent information, addressing issues related to repetition and irrelevant details in the retrieved content.

Task Adaptable Module

This module makes RAG adaptable to various downstream tasks allowing the development of task-specific end-to-end retrievers with minimal examples, demonstrating flexibility in handling different tasks.

Acknowledgements

Retrieval Augmented Generation continues to be a pivotal approach for any Generative AI led application and it is only going to grow. There are several individuals and organisations that have provided learning resources and made understanding RAG fun.

I'd like to thank -

- My team at [Yarnit.app](#) for taking a bet on RAG and helping me explore and execute RAG pipelines for content generation
- **Andrew Ng** and the good folks at [deeplearning.ai](#) for their short courses allowing everyone access to generative AI
- **OpenAI** and **HuggingFace** for all that they do
- **Harrison Chase** and all the folks at [LangChain](#) for not only building the framework but also making it easy to execute
- **Jerry Liu** and others at [Llamaindex](#) for their perspectives and tutorials on RAG
- [TruEra](#) for demystifying observability and the tech stack for LLMOps
- [PineCone](#) for their amazing documentation and the learning center
- The team at [Exploding Gradients](#) for creating [Ragas](#) and explaining RAG evaluation in detail
- [TruLens](#) for their triad of RAG evaluations
- [Aman Chadha](#) for his curation of all thing AI, ML and Data Science
- Above all, to my **colleagues and friends**, who endeavour to learn, discover and apply technology everyday in their effort to make the world a better place.

With lots of love,

Abhinav



I talk about :

#AI #MachineLearning #DataScience
#GenerativeAI #Analytics #LLMs
#Technology #RAG #EthicalAI

let's connect...



[DOWNLOAD FREE
EBOOK](#)



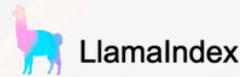
[DOWNLOAD FREE
EBOOK](#)

Resources

Official Documentations



[Python Documentation](#)



[Python Documentation](#)



[Learning Center](#)



[Documentation](#)



[Documentation](#)

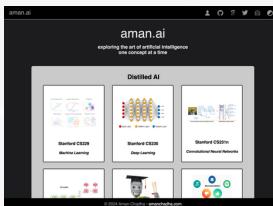


[Hugging Face](#)



[OpenAI](#)

Thought Leaders and Influencers



[Aman Chadha's Blog](#)



[Lillian Weng's Log](#)



[Leonie Monigatti's Blogroll](#)



[Chip Huyen Blogs](#)

Research Papers



Retrieval-Augmented Generation for Large Language Models: A Survey
(Gao, et al, 2023)



Retrieval-Augmented Multimodal Language Modeling
(Yasunaga, et al, 2023)



KG-Augmented Language Models for Knowledge-Grounded Dialogue
(Kang, et al, 2023)

Learning Resources and Tutorials



[Short 1-hour Courses](#)



[Python Cookbook](#)



[Tutorials & Webinars](#)

Hello!

I'm Abhinav...

A data science and AI professional with over 15 years in the industry. Passionate about AI advancements, I constantly explore emerging technologies to push the boundaries and create positive impacts in the world. Let's build the future, together!



Please share your feedback on these notes with me



LinkedIn

Github

Medium

Insta

email

X

Linktree

Gumroad

Talk to me

Book a meeting

Virtual Coffee
Ask me anything
Resume review (DS, AI, ML)
AI ML Strategy Consultation



topmate.io/abhinav_kimothi

Checkout Yarnit Magic


Yarnit

5-in-1 Generative AI Powered
Content Marketing Application

www.yarnit.app

Newsletter



Vital Vector
by Abhinav Kimothi
Read and subscribe at
vitalvector.substack.com

\$\$ Contribute \$\$



Buy me a coffee

Subscribe

Don't miss a post!

Get an email notification whenever I publish!

Follow on LinkedIn

