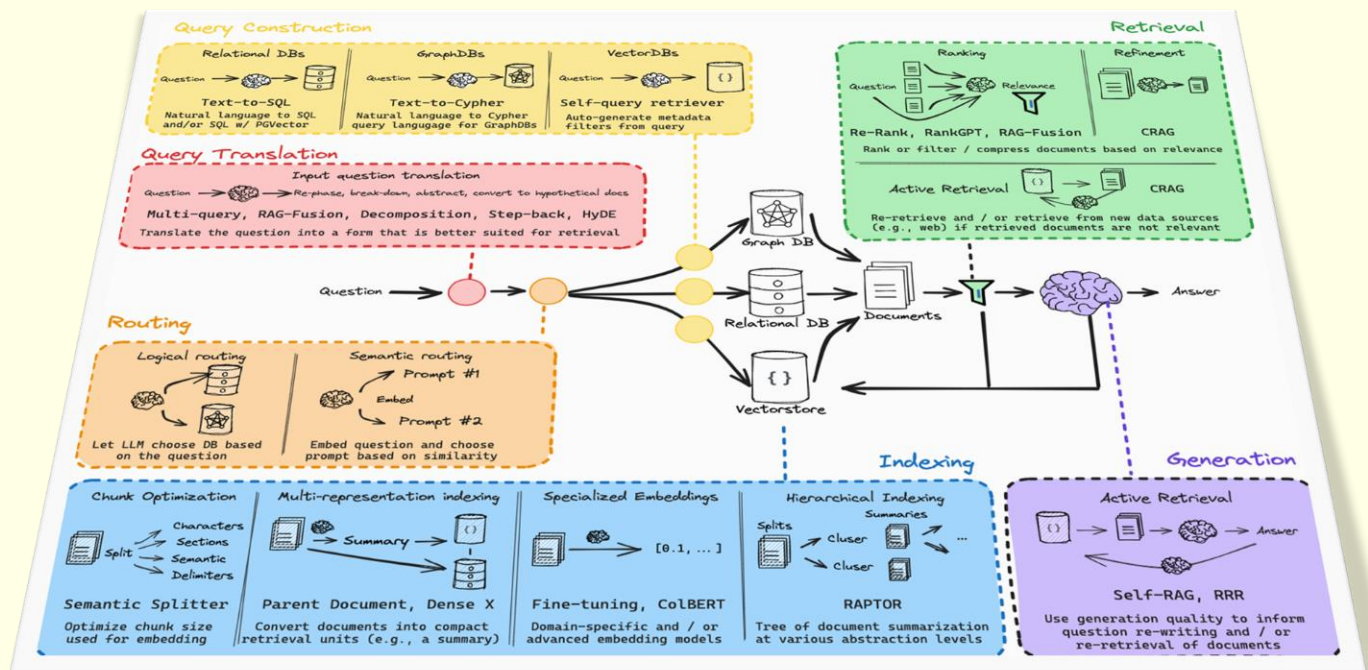# RAG Tutorial

## (Retrieval Augmented Generation)

### Indexing and Loading Blog Post Content Using DocumentLoaders



# Day 3 of 7

# <u>Table of Contents</u>

# 1. Introduction

In the previous sessions, we set up our environment and previewed the construction of a question-answering app.

Now, we need to load the content we want to index and retrieve from. For this, we'll use DocumentLoaders, which are designed to fetch data from various sources and return it as a list of Documents.

## *2. Loading Blog Post Contents*

We'll DocumentLoaders simplify the task of fetching and parsing data from different sources.

Each Document consists of page_content (a string) and metadata (a dictionary).

In this tutorial, we'll use the WebBaseLoader to load and parse HTML content from a specified URL.

# 3. Customizing HTML Parsing

The WebBaseLoader utilizes urllib to fetch HTML content and BeautifulSoup to parse it.

We can customize the parsing process by passing parameters to the BeautifulSoup parser through bs_kwargs.

For our purpose, we will filter and keep only HTML tags with the classes "post-content", "post-title", and "post-header"

# *4. Sample Code and Explanation*

In this Google Collab Notebook, we will demonstrates how to extract specific content from a website using the `langchain_community` and `bs4` (Beautiful Soup) libraries.

First, it defines a `bs4_strainer` using `SoupStrainer` to target only HTML elements with classes "post-title", "post-header", and "post-content".

Then, it initializes a `WebBaseLoader` with the target website URL and the strainer. When `loader.load()` is called, it fetches the website content and parses it, keeping only the specified elements.

Finally, the code prints the length of the extracted content and displays the first 500 characters as a preview.

This technique is useful for extracting relevant information from websites while discarding unnecessary elements.

# Link of collab Notebook

# *5. Conclusion*

In this tutorial, we explored how to use the *WebBaseLoader* to load and parse HTML content from a web URL.

We customized the parsing process to filter and retain only relevant sections of the HTML using *BeautifulSoup*.

This forms the basis for our content indexing, which is crucial for building our question-answering app.

Stay tuned for Day 4, where we will delve deeper into the indexing process and prepare our data for efficient retrieval.

Stay tuned for the next tutorial in this series, where we will dive deeper into optimizing and scaling your RAG applications.

ANSHUMAN JHA