# Symantic Spotter using LlamaIndex

Author: Rakesh Rangadal

Date: 01-10-2025

**Project Objective**

Build a project in the insurance domain, similar to the project you saw in the 'Retrieval Augmented Generation' session. The goal of the project will be to build a robust generative search system capable of effectively and accurately answering questions from various policy documents. You may use LangChain or LlamaIndex to build the generative search application

**Solution Strategy**

Build a solution which should solve the following requirements using LlamaIndex:

- Users would responses from insurance policy knowledge base.
- If user want to perform a query system must be able to response to query accurately.

**Goal**

 Solving the above two requirements well in and would ensure that the accuracy of the overall model is good.

*Data Used*

HDFC various Insuracne policy documetns sotred in single folder

*Tools used* LlamaIndex, GhatGPT has been used due to its powerful query engine, fast data processing using data loaders and directory readers as well as easier and faster implementation using fewer lines of code.

*Tools used*

 LlamaIndex, GhatGPT, disc cache  has been used due to its powerful query engine, fast data processing using data loaders and directory readers as well as easier and faster implementation using fewer lines of code.

# Why LlamaIndex ?

LlamaIndex is an innovative data framework specially designed to support LLM-based RAG framework application development. It offers an advanced framework that empowers developers to integrate diverse data sources with large language models.
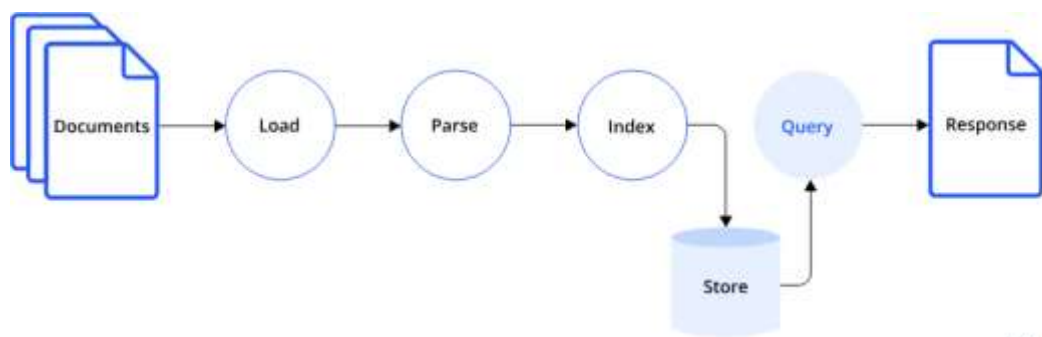
LlamaIndex includes a variety of file formats, such as PDFs and PowerPoints, as well as applications like Notion and Slack and even databases like Postgres and MongoDB.

The framework brings an array of connectors that assist in data ingestion, facilitating a seamless interaction with LLMs. Moreover, LlamaIndex boasts an efficient data retrieval and query interface.

LlamaIndex enables developers to input any LLM prompt and, in return, receive an output that is both context-rich and knowledge-augmementation.
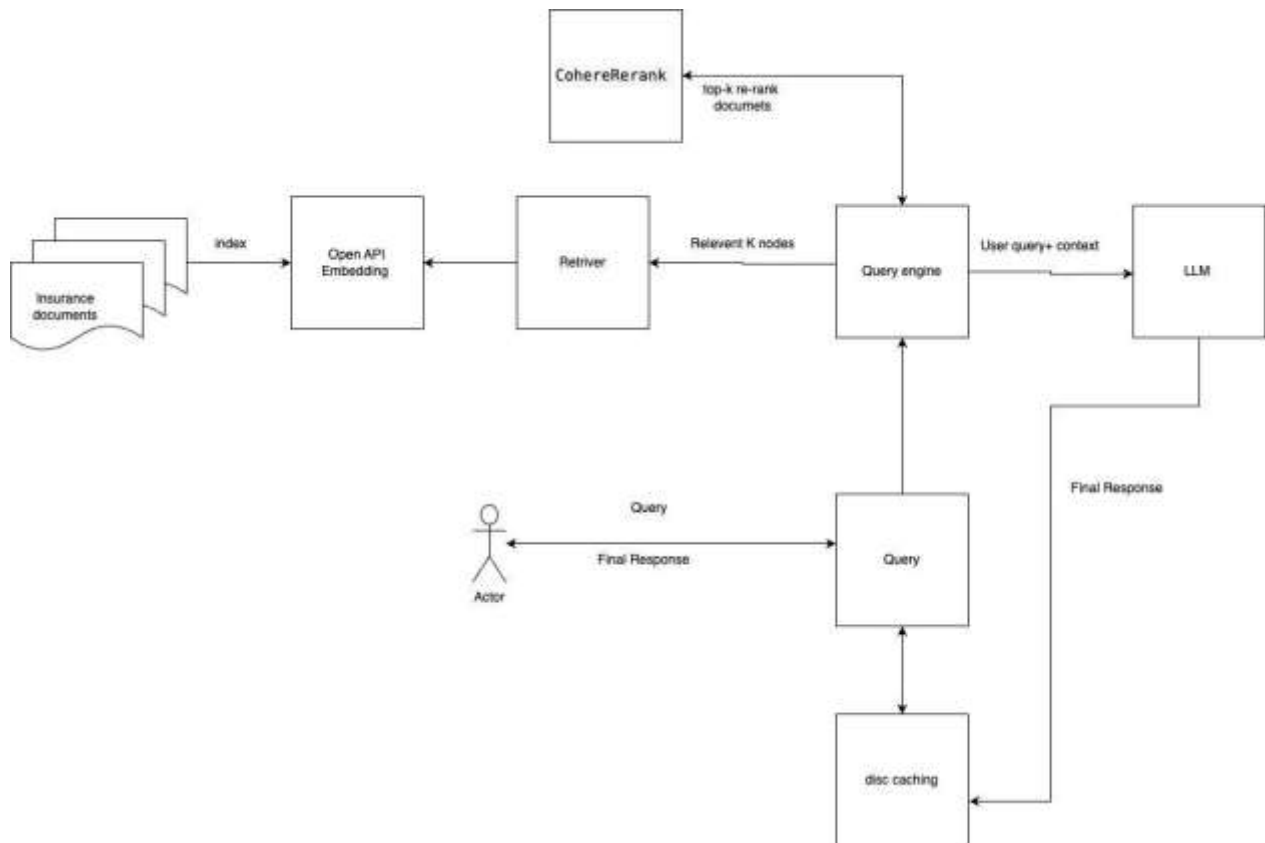
## Key Feature of LlamaIndex:

- Data connectors allow ingestion from various data sources and formats.
- It can synthesize data from multiple documents or heterogeneous data sources.
- It provides numerous integrations with vector stores, ChatGPT plugins, tracing tools, LangChain, and more.



LeewayHertz

High Level Design of  Semantic Spotter :

CohereRerank

top-k re-rank documets

Insurance documents

index

Open API Embedding

Retriver

Relevent K nodes

Query engine

User query+ context

LLM

Actor

Query

Final Response

Query

Final Response

disc caching

Final Response

Architecture Descriptions:

1. *Documents*: We will be using list of HDFC insurance documents provides inside a single folder.
2. *Open API embedding*: We ars using OpenAPI embedding as Vector DB for indexing insureance documetns in the form of embedding.
3. *Query Engine*: We are using Query Engine Module of Llammaindex for perfomring synantic Search. Query Engine will use internally Retriver and Cohere Rerank to retrrive top-k relvent nodes from embedding.
4. *LLM* : top k-documets along with usewr query will be passed to LLM to generate the accurate response.We are using chatGPT LLM.
5. *Caching*:" Caching is being used to improve the read operation. Recent similar search will be store in Caching and user query frist will be served from Cahcing. If user query not found in cahce then query will be forwarded to query engine and then LLM to geneate
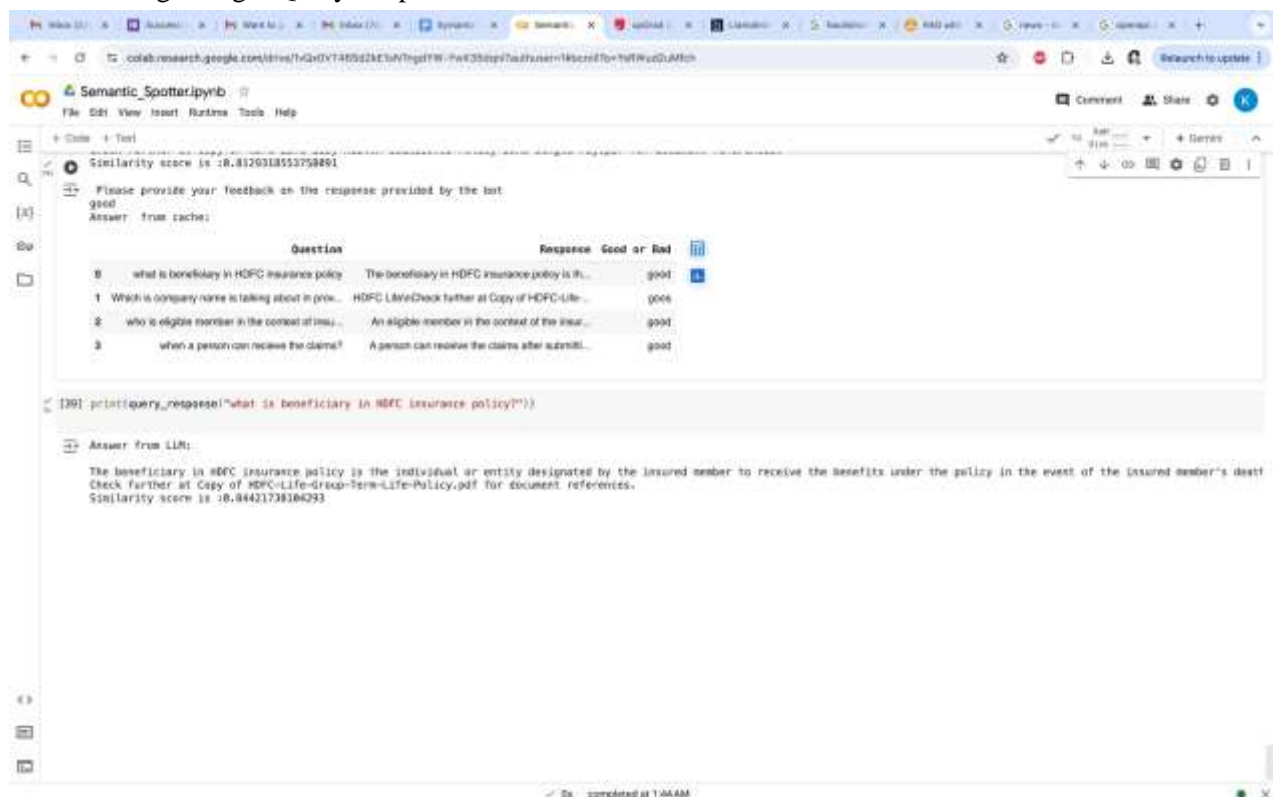
the resposne. user query and generated resposne will be cached in in cache and will be served from there based ttl.
6. *Meta data*:Along with Resposne we are also returning docs referece and similary score to improve the user confidence towards the implemented RAG system.
7. *Cohere-Rerank*: Is being used to rerank the query based on semantic score.

## Generative Search Response from Insurance documents:

We have attached custom query generative search results.

1. Using a single Query Respone:



Multiple Query Resposne:

Code:
Jupyter Notebook developed for Semantic Spotter.


**Chalenges Faced:**

We tried to use GPTCache with for caching system,but due compatiability issues, we couldn't interate it.

**Alternative Solution:**

We are usign disc caching for alternative to GPTCache.

## Future of work

1. We can further improve the solution by using feedback mechanism usingc embedding the response along with feedback to vector database.
2. We can also integrate RAGAS/DeepEval Framework to improve the symnatic search and generated response.