

In the Name of Allah

Candidate: Fatemeh Vakili

Job Position: Senior Python Developer

This report outlines the implementation of the RAG System with Vector Database and the rationale behind selecting each component for this project. The purpose of the issue is summarized below:

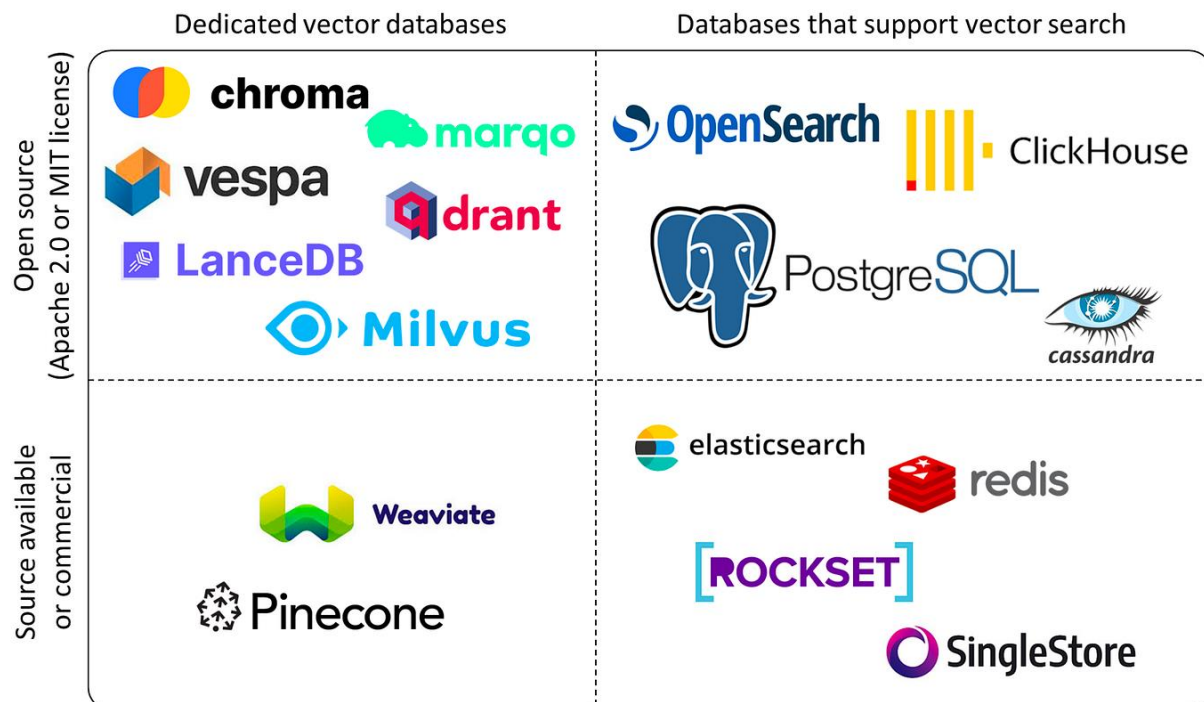
Develop a Retrieve-and-Generate (RAG) system that leverages a vector database to store and retrieve information. The system should expose an API built with FastAPI to answer questions based on the embedded data.

Note: The project requirements are a step-by-step process! The problem requirements are written in purple, with explanations following in black.

Requirements

1. Vector Database Integration:

- Implement a vector database (excluding ChromaDb) to store information in a structured format that allows for efficient similarity search.
- Justify the choice of the database considering factors such as performance, scalability, and community support.



The image above illustrates various vector databases, some of which are open source while others are not. Specific databases offer additional vector features, such as Elastic Search and Mongo, making them highly popular. This report is designed to present two top vector databases based on

specific criteria. Critical considerations for selecting a database include its performance, scalability, and community support.

Pinecone:

Pinecone is a platform for managing vector databases designed to tackle the challenges of working with high-dimensional data. It offers advanced search and indexing capabilities, empowering data scientists and engineers to efficiently create and implement large-scale machine learning applications. Pinecone leverages Faiss for efficient similarity searches and boasts the following key features:

- Fully managed service
- High scalability
- Real-time data ingestion
- Low-latency search
- Integration with LangChain

In addition to its cloud-native infrastructure and seamless API usage, Pinecone provides essential functionalities such as duplicate detection, rank tracking, data search, classification, and deduplication. This enables users to concentrate on advancing their AI solutions without the burden of managing infrastructure.

Qdrant:

Qdrant is a high-performance vector database and similarity search engine designed to manage high-dimensional vector embeddings (Points) with metadata (payloads). It is built in Rust to ensure speed and reliability, even under heavy loads, and provides a user-friendly API for storing, searching, and managing vectors. Qdrant is tailored for tasks involving neural network or semantic-based matching, faceted searches, and recommendations, making it ideal for developers needing robust vector search capabilities. Key features include:

Versatile API: It supports Python, TypeScript/JavaScript, Rust, and Go API clients, along with OpenAPI v3 specs.

Speed and Precision: It utilizes the HNSW algorithm with various distance metrics (Cosine, Dot, Euclidean) for fast and accurate searches.

- **Advanced Filtering:** Supports filtering results based on vector payloads, which are helpful for diverse data types like string matching, numerical ranges, and geo-locations.
- **Scalability:** Cloud-native design with horizontal scaling options and a free tier for exploration.
- **Efficiency:** Built-in Rust for optimal resource use and dynamic query planning.

Based on the presented cases, both databases are considered suitable options. We have chosen to utilize the Pinecone database because of its popularity and compatibility with our database.

2. Data Modeling:

- Design data models that effectively represent the information to be stored in the vector database.
- Explain the rationale behind the chosen schema and how it facilitates the RAG system's operations.

Vector databases store data in high-dimensional vectors, representing features or attributes mathematically. Depending on the complexity of the data, these vectors can have varying dimensions, ranging from tens to thousands. Unlike vector libraries integrated into broader databases to enable similarity search, vector databases are specifically designed for efficiently managing dense vectors and supporting advanced similarity search.

Vector similarity search is crucial for vector databases because it efficiently identifies and retrieves the most relevant vectors from large datasets. This capability is essential for making quick and accurate data-driven decisions, reducing the time and computational resources needed for querying large datasets. The evolution of vector similarity search has been driven by the necessity to efficiently analyze and interpret large datasets, evolving alongside relational databases since the mid-20th century.

To summarize, as data becomes more abundant and diverse, methodologies for vector similarity search are indispensable for quick and accurate data retrieval, powering various applications in artificial intelligence, machine learning, and information retrieval. Understanding these methods helps grasp how data-driven technologies operate effectively.

Vector Similarity Search Algorithms

1. Dot Product:

The Dot Product is a fundamental similarity measure utilized in machine learning, data mining, and statistics to evaluate the similarity between two vectors. Its significance extends to cosine similarity and underpins algorithms in search engines, recommendation systems, and various data-driven applications.

2. Cosine Similarity:

Cosine Similarity serves as a widely employed metric for gauging similarity between vectors. Commonly used in information retrieval, text mining, and machine learning, it gauges the cosine of the angle between two vectors, offering insights into their orientation within a multidimensional space.

Applications in Similarity Search:

- Document Similarity: Primarily employed in natural language processing to evaluate similarities between texts, assisting in document retrieval and clustering.

- Recommendation Systems: Cosine Similarity is deployed in collaborative filtering to produce recommendations by comparing user or item profiles.

3. Euclidean Distance (L2 Distance):

Euclidean Distance is a prevalent measure of similarity used in data mining, machine learning, and statistics. It computes the straight-line distance between two points in Euclidean space and offers a straightforward measure of vector dissimilarity.

Top companies like OpenAI, Google, Cohere, and Anthropic currently lead the global market for large language models (LLMs), focusing on solving Natural Language Processing (NLP) tasks. There is also intense competition to offer the best multilingual text embedding services. Here's a brief overview:

- OpenAI provides a closed-source API for multilingual text embeddings. Their latest model, text-embedding-3-large, was released on January 25, 2024. It supports 256, 1024, and 3072 dimensions, with 3072 as the default.

- Cohere, This company offers a closed-source API for multilingual text embeddings. Its latest model, embed-multilingual-v3.0, was released on November 2, 2023, and returns embeddings with 1024 dimensions.

- Google provides a closed-source API for multilingual text embeddings. Their latest model, text-multilingual-embedding-preview-0409, was made available for preview on April 2, 2024, and returns embeddings with 768 dimensions.

- Microsoft offers an open-source, multilingual embedding model called Multilingual E5. The initial release included three models (small, base, and large), with the instructional version released in early 2024.

Each provider aims to deliver the most effective multilingual embedding service, catering to the diverse needs of NLP applications worldwide.

Retrieval Augmented Generation (RAG) is a technique that improves pre-trained large language models (LLMs) like GPT-3 or GPT-4 by integrating them with external data sources. This combination harnesses the generative abilities of LLMs with precise data retrieval mechanisms to provide nuanced and accurate responses. RAG is particularly beneficial in situations that require specific and up-to-date information, addressing the limitations of LLMs such as:

- Lack of specific information: Traditional LLMs may give generic responses due to constraints in their training data.

- Hallucinations: LLMs may produce incorrect responses based on imagined facts.

- Generic responses: LLMs may not offer context-specific answers, which is problematic in personalized applications such as customer support.

How RAG Works

1. Data Collection: Gathering necessary data, such as user manuals and product databases.

2. Data Chunking breaks down large documents into smaller, focused sections to improve retrieval efficiency and relevance.

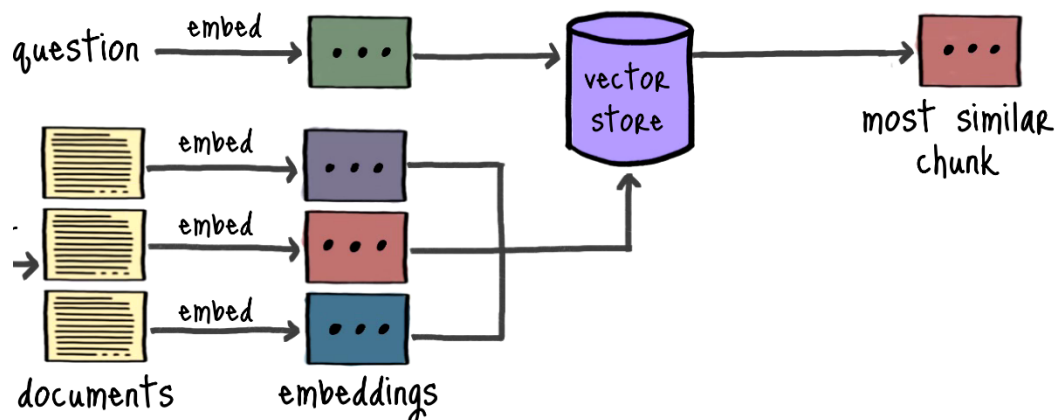
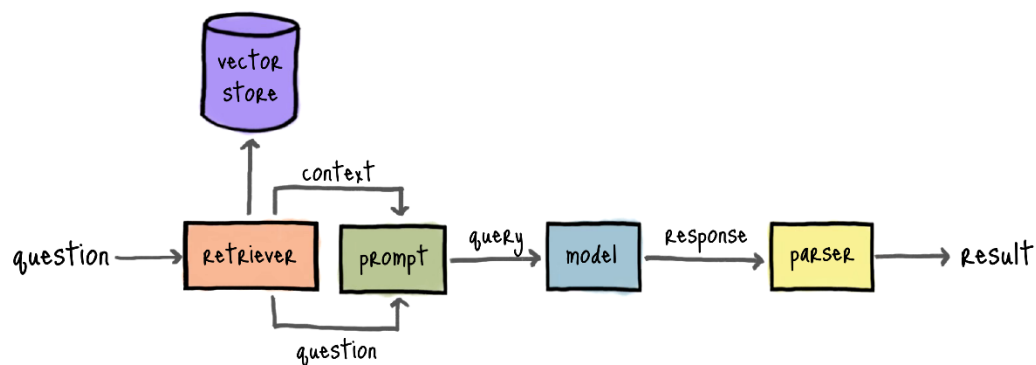
3. Document Embeddings: Converting these data chunks into vector representations (embeddings) to capture their semantic meanings.

4. Handling User Queries: Converting user queries into embeddings and comparing them with document embeddings to find the most relevant information.

5. Generating Responses: Feeding the retrieved information and the user query into an LLM to generate a coherent and accurate response.

This approach, facilitated by frameworks like LlamaIndex and OpenAI Models, ensures that LLMs can efficiently access and utilize specific data. It is instrumental in applications such as customer support chatbots, where accurate and personalized information is essential.

The two photos below display the LangChain-based chain and the system. Once similar chunks are identified, the most similar answers are shown.



3. API Development:

- Use FastAPI to create endpoints that allow users to submit questions and receive answers.
- Ensure the API is well-documented and adheres to RESTful principles.

4. RAG System Logic:

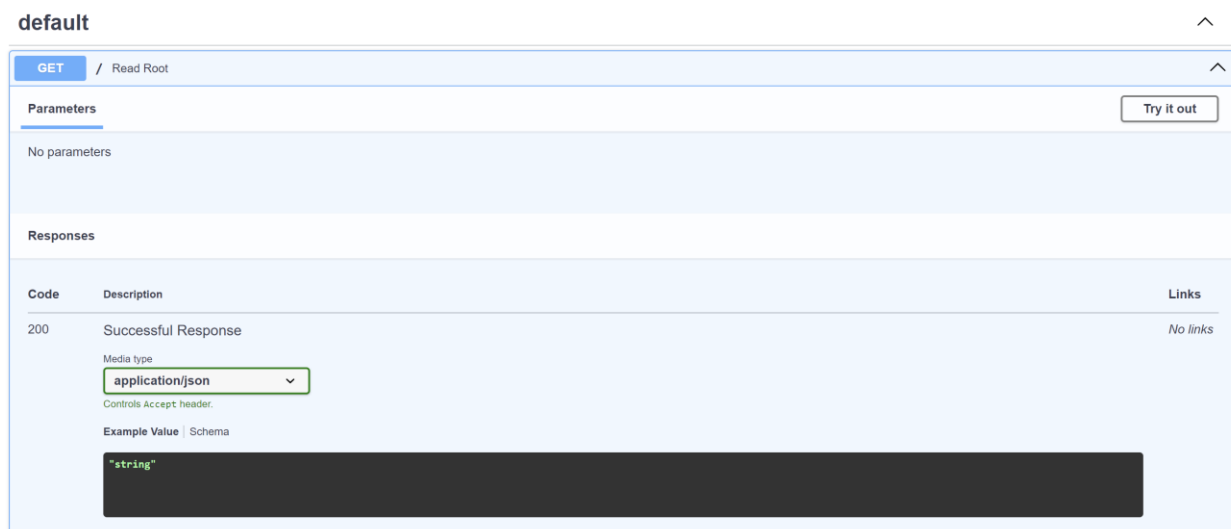
- Develop the logic for the RAG system to retrieve relevant information from the vector database and generate answers.
- Incorporate state-of-the-art NLP techniques to enhance the quality of the generated responses.

The principles of REST APIs encompass the utilization of resources to represent data, including entities such as users, products, or orders, and the employment of HTTP methods to execute CRUD operations on these resources. RESTful APIs are designed to adhere to principles such as statelessness, a uniform interface, and resource-based interactions. They rely on URIs (Uniform Resource Identifiers) to identify resources and offer transparent access and manipulation structures for data. The primary HTTP methods utilized are as follows:

- GET: This method retrieves a specific resource or a collection of resources.
- POST: This method is used to create a new resource.
- PUT: This method is used to update an existing resource.
- DELETE: This method is used to delete a specific resource.
- PATCH: This method is used to partially update an existing resource.

This project has code versions with and without FastAPI. The FastAPI-based file provides details about **Swagger UI or ReDoc.**

This section uses the "get" command to show the address <http://127.0.0.1:8000/> on the local host. It's a command used for reading data.



You can upload text, images, Jason, or PDF data into the system. For this reason, the post command is used to send files.

The screenshot shows a REST client interface for a POST request to `/uploader/` with the description "Create Upload File". The "Parameters" tab is active, showing "No parameters". The "Request body" tab is also active, showing a required field "file" of type "string(\$binary)". A "Choose File" button is present, and the file "YOLO.pdf" is selected. The "multipart/form-data" content type is chosen from a dropdown. At the bottom, there are "Execute" and "Clear" buttons.

The desired file has been uploaded to the server in the following image:

The screenshot shows the response of the POST request to `/uploader/`. The "Curl" tab shows the command used. The "Request URL" is `http://127.0.0.1:8000/uploader/`. The "Server response" tab shows a 200 status code. The "Response body" is a JSON object: `{ "filename": "YOLO.pdf" }`. The "Response headers" are: `content-length: 23`, `content-type: application/json`, `date: Mon, 08 Jul 2024 00:04:51 GMT`, and `server: uvicorn`. The "Responses" tab shows a "Successful Response" with a 200 status code. The "Media type" is set to "application/json".

In this section, you can use the post command to use the system based on the uploaded file.

The screenshot shows a REST client interface for a POST request to `/ask-question/` with the description "Ask Question Endpoint". The "Parameters" tab is active, showing a required field "question" of type "string(query)". The query "What is YOLO?" is entered. At the bottom, there are "Execute" and "Clear" buttons.

In this picture, a question that is not using Fast API is posed to the system. As shown, the system provides a corresponding response.

Responses

Curl

```
curl -X 'POST' \
  'http://127.0.0.1:8000/ask-question/?question=what%20is%20YOLO%3F' \
  -H 'accept: application/json' \
  -d ''
```

Request URL

```
http://127.0.0.1:8000/ask-question/?question=what%20is%20YOLO%3F
```

Server response

Code	Details
200	<p>Response body</p> <pre>{ "question": "What is YOLO?", "answer": "Answer: YOLO stands for \"You Only Look Once,\" which is a new approach to object detection presented in the document. It reframes object detection as a regression problem + spatially separated bounding boxes and associated class probabilities, allowing for real-time processing of images with high speed and accuracy." }</pre> <p>Response headers</p> <pre>content-length: 359 content-type: application/json date: Mon, 08 Jul 2024 00:07:01 GMT server: uvicorn</pre>

Please note that only the backend has been reviewed for this project. If we intend to utilize the front end, it's essential to integrate the request command into the backend.

Also, the advantages of using FastAPI can be listed as follows:

- Data Validation
- Auto Documentation
- Auto-Completion & Code Suggestion

NLP techniques are helpful when working with textual data. For instance, Python's NLTK library or Hugging Face models can serve as effective transformers!

Retrieval-augmented generation (RAG) systems meld retrieval-based and generative models to deliver more accurate and contextually relevant responses. Improving RAG systems involves enhancing both retrieval and generation through various NLP techniques:

1. Improving Retrieval:

- Utilize dense vector embeddings (e.g., BERT) for better contextual understanding.
- Include hard negatives during training to improve model discernment.
- Expand queries with synonyms or related terms to boost recall.
- Incorporate more context from the original query for better document retrieval.
- Fine-tune retrieval models on domain-specific data for improved relevance.

2. Enhancing Generation:

- Fine-tune pre-trained generative models (e.g., GPT-3, T5) on domain-specific datasets.
- Inject external knowledge sources into the generative model for more informative responses.
- Use constraints to ensure factually correct and contextually appropriate outputs.

- Employ post-processing techniques like reranking and fact-checking to enhance output quality.

3. Hybrid and Ensemble Techniques:

- Combine dense and sparse retrieval models to leverage their strengths.
- Use hybrid answer scoring from multiple models to select the most accurate responses.
- Implement iterative cycles where the generator refines the retriever's inputs and vice versa.

4. Evaluation and Feedback:

- Incorporate human feedback to fine-tune the system continuously.
- Utilize automatic evaluation metrics (e.g., BLEU, ROUGE, BERTScore) for quality assessment.
- Conduct A/B testing to determine the best configurations empirically.

5. Context Management:

- Break down complex documents into chunks and use hierarchical retrieval to fetch the most relevant sections.
- Optimize context window size for generative models to balance comprehensiveness and relevance.

6. Adversarial Training:

- Train with adversarial examples to enhance robustness against misleading inputs.

7. Cross-Modal and Multimodal Retrieval:

- Incorporate information from multiple modalities (text, images, tables) for complex queries.

Additional techniques include:

- Data Indexing Optimizations: Improve indexing methods for efficient retrieval.
- Hybrid Search Models: Combine different retrieval approaches for better results.
- Fine-Tuning Embedding Models: Customize embeddings for domain-specific tasks.
- Response Summarization: Summarize responses for concise and relevant answers.
- Re-ranking and Filtering: Re-ranking and filtering to improve the quality of retrieved documents.

Implementing and fine-tuning these techniques can significantly enhance RAG systems' performance, accuracy, and reliability for real-world applications.

Another choice for those who prefer not to write code is to use Longflow!