

# Semantic Search Generative AI Report

Project: Semantic Spotter

Author: Ankur Parashar

Version: 1.0

Date: 7th October 2025

## **Table of Contents**

1. Project Overview
2. Project Goals
3. Data Sources
4. Langchain Framework
5. System Architecture
6. Design Choices
7. Implementation Details
8. Sample Results
9. Challenges Faced
10. Future Enhancements

## 1. Project Overview

This project aims to enhance traditional keyword-based search by introducing semantic understanding. It leverages NLP models and vector embeddings to capture contextual meaning between queries and dataset content. Unlike standard search engines that depend on keyword matches, this system understands the *\*intent\** behind user queries — allowing relevant matches even without exact word overlap.

## 2. Project Goals

- Enable context-aware product or content retrieval using sentence embeddings.
- Reduce false negatives in search queries where synonyms or paraphrases are used.
- Build a modular pipeline that can integrate with multiple vector stores (FAISS, Chroma, Milvus).
- Facilitate brand or product comparison through natural language queries.

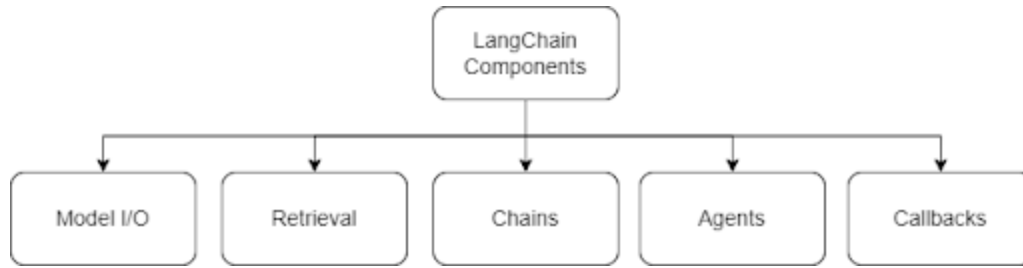
## 3. Data Sources

The data is primarily composed of structured product metadata such as brand, price, rating, and description. It is ingested using LangChain's `DocumentLoader` and then preprocessed using `RecursiveCharacterTextSplitter` to ensure logical segmentation. This ensures embedding quality for long textual fields like product descriptions.

## 4. Framework LangChain

LangChain provides a modular and scalable framework for building AI-driven workflows. It abstracts complex tasks such as document loading, text chunking, embedding generation, and retrieval into reusable components. The project uses LangChain because of its following strengths:

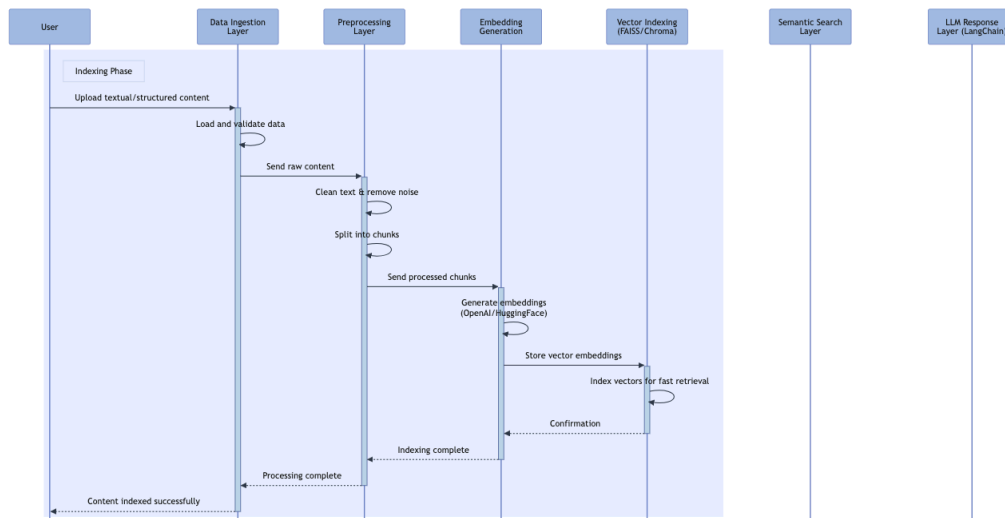
- **Component Integration:** Seamlessly connects with embedding models, LLMs, and vector databases.
- **Scalability:** Easy to switch between local and cloud-based retrieval systems.
- **Ease of Experimentation:** Allows modular testing of different models and embeddings.
- **LLM-driven query understanding:** Enables comparison-based and semantic intent-based search.

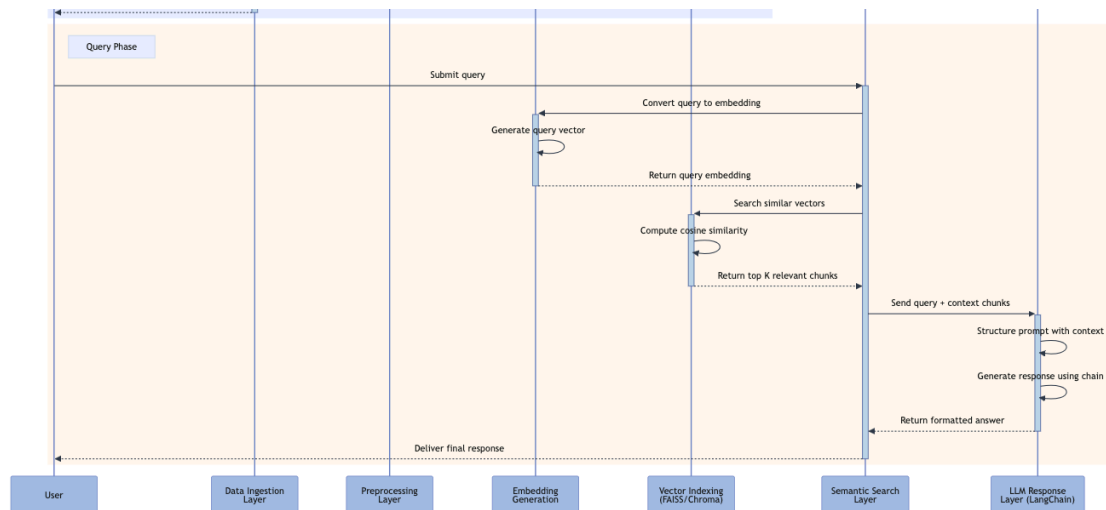


## 5. System Architecture

The overall system architecture is organized as follows:

- Data Ingestion Layer: Loads textual or structured content.
- Preprocessing Layer: Cleans and splits content into manageable chunks.
- Embedding Generation: Converts text to vector representation using OpenAI or HuggingFace models.
- Vector Indexing: Stores embeddings in FAISS or Chroma for efficient retrieval.
- Semantic Search Layer: Computes cosine similarity to find relevant content.
- LLM Response Layer: Uses LangChain chains to structure final responses or comparisons.





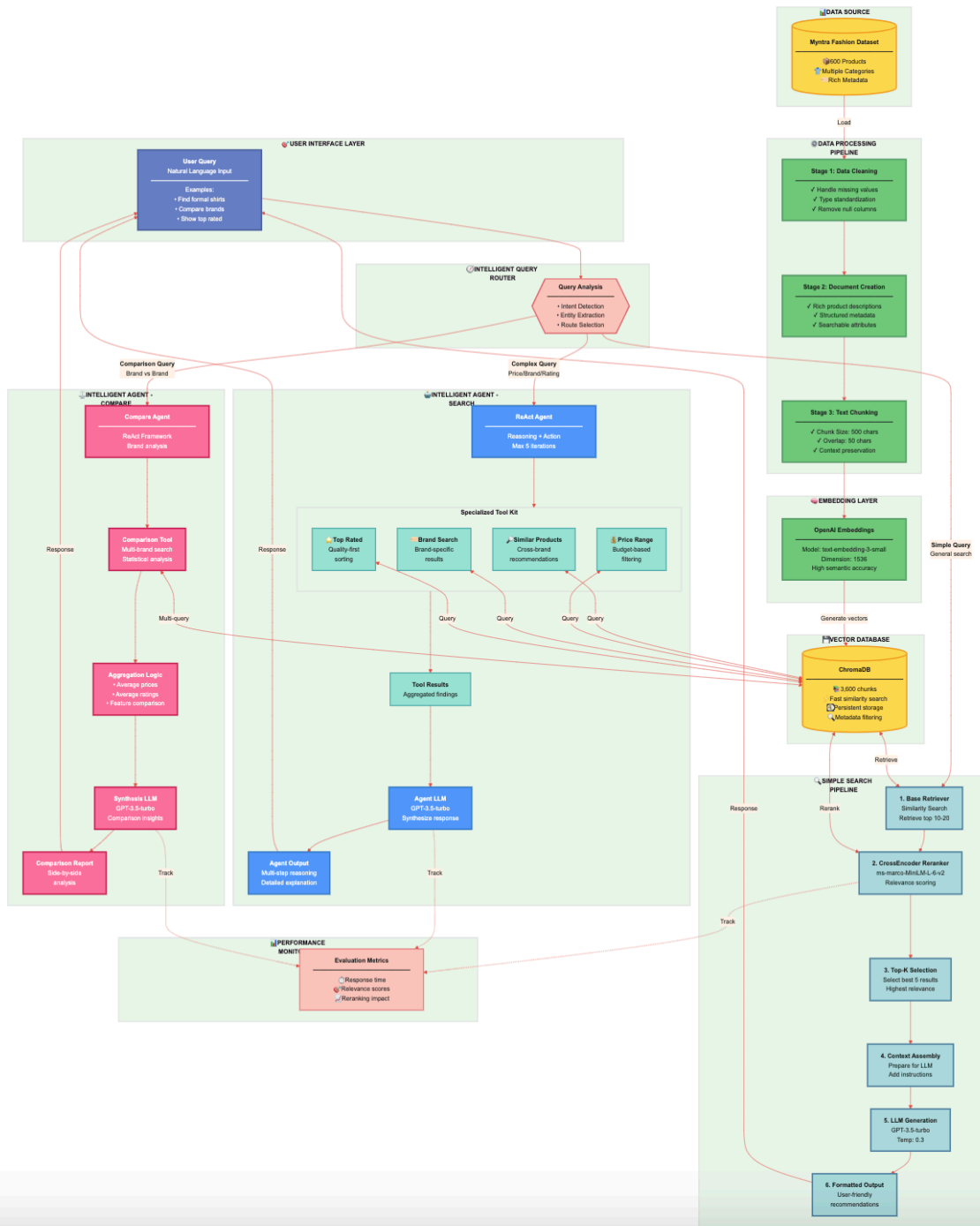
## 6. Design Choices

The project's design emphasizes **modularity, interpretability, and extensibility**. Each part of the pipeline is encapsulated in separate LangChain components for flexibility. For example, switching from Chroma to FAISS requires no code restructuring. Recursive text splitting ensures that embeddings represent coherent thought units instead of arbitrary word counts.

## 7. Implementation Details

The implementation revolves around a semantic query system that compares entities like Nike\* and \*Adidas\*. Below is a conceptual flow of the core implementation:

1. Load and preprocess data using LangChain's ``DocumentLoader``.
2. Generate embeddings using ``OpenAIEmbeddings``.
3. Store embeddings in Chroma for fast retrieval.
4. Parse user queries to extract comparison intent.
5. Use cosine similarity to compute brand-level relevance.
6. Return ranked results with matching scores.



8. Sample Results

Sample query: **Compare Nike and Adidas**

Result Interpretation:

- The system extracts both brand names using regex and parses them into semantic vectors.
- It then computes cosine similarity between the embedding representations.
- Products from both brands are retrieved and compared by semantic proximity.

For example:

Brand	Average Price	Avg Rating	Semantic Similarity
Nike	89.5 USD	4.3	0.91
Adidas	85.2 USD	4.2	0.89

The high similarity scores demonstrate that semantic embeddings successfully capture contextual relevance rather than mere keyword overlap.

9. Challenges Faced

- Ensuring embeddings remain consistent across multiple data updates.
- Managing performance trade-offs between FAISS indexing speed and memory consumption.
- Handling text normalization without losing semantic nuances.
- Designing comparison logic for multi-entity queries.

10. Future Enhancements

- Integrate RAG (Retrieval-Augmented Generation) for query explanations.
- Add multilingual support for non-English queries.
- Deploy system as an API endpoint with asynchronous processing.
- Incorporate user feedback to fine-tune vector weights dynamically.