# Semantic Spotter Project Documentation

## 1. Problem Statement

The primary objective of the **Semantic Spotter Project** is to enhance the product search experience for users looking to find fashion items in the **Myntra Fashion Database**. Traditional search engines rely heavily on keyword matching, which can lead to irrelevant or incomplete search results, especially when the user query is phrased differently or uses synonyms. Additionally, current search systems often lack the ability to generate meaningful descriptions and explanations alongside search results, which could improve the decision-making process for shoppers.

The challenge is to build a more intuitive, semantic-based search system that allows users to find relevant products with greater accuracy, using natural language queries, while also generating helpful, context-rich responses that are easy to understand and engaging.

## 2. Goals

The key goals of the **Semantic Spotter Project** are as follows:

1. **Improved Search Accuracy**: Implement a search system that can understand the meaning behind a user's query and retrieve semantically relevant products, not just keyword matches.
2. **Contextual Generation**: Create a system that generates not just a list of products, but also a meaningful response that includes product descriptions, features, and recommendations based on the user's query.
3. **Enhanced User Experience**: Enable users to easily search and understand product details using natural language, rather than forcing them to use strict filtering criteria or specific keywords.
4. **Scalability and Speed**: Ensure that the solution works efficiently even with a large dataset of products, providing fast retrieval times while maintaining the quality of results.

# 3. Solution

To address the problem, the **Semantic Spotter Project** combines **Retrieval-Augmented Generation (RAG)** with **LangChain**, **ChromaDB**, and **OpenAI embeddings**. Here's how the solution works:

## Solution Architecture

1. **Natural Language Query Input**: Users enter a query in natural language, such as "Summer dresses in red and white" or "Comfortable blue shoes."
2. **Query Embedding and Search**: The input query is first transformed into an embedding using **OpenAI's embeddings model**. This embedding captures the semantic meaning of the query, making it possible to compare it with product embeddings in the vector store.
3. **Vector Database (ChromaDB)**: The product data from Myntra is pre-processed, and textual information such as descriptions and features are converted into embeddings. These embeddings are stored in **ChromaDB**, a high-performance vector database optimized for fast similarity searches.
4. **Similarity Search**: When a user query is received, it is transformed into an embedding and compared to the embeddings in **ChromaDB** using similarity measures (e.g., cosine similarity). The system retrieves the top N products that are most semantically relevant to the query.
5. **Contextual Response Generation**: After retrieving the top matching products, an **OpenAI language model** is used to generate a contextual response. This step combines the raw product data with natural language generation to offer rich, detailed product descriptions and recommendations, thus ensuring a user-friendly and informative search experience.
6. **Final Response to User**: The final response is delivered to the user, which includes the relevant product details (e.g., name, description, price, etc.) and possibly some additional context or suggestions, such as why these products match the query or how they can be used.

# 4. The Model Used

The **Semantic Spotter Project** leverages the following models and technologies:

1. **OpenAI's Embeddings Model**: This model is used to convert the product descriptions and the user query into high-dimensional vector embeddings. The embeddings capture the semantic meaning of the text, allowing for efficient similarity-based search.
2. **OpenAI GPT Models**: For the generation phase, we use an OpenAI GPT model to generate meaningful text about the retrieved products. This model is responsible for crafting context-rich responses that are informative and engaging.
3. **LangChain**: LangChain is used as the framework that facilitates the integration of the **OpenAI embeddings**, **ChromaDB**, and **retrieval-augmented generation** (RAG). LangChain allows the orchestration of the flow, connecting the various components (embedding generation, vector search, and text generation).
4. **ChromaDB**: ChromaDB is a vector database used for efficient storage and retrieval of product embeddings. It provides fast similarity search capabilities, ensuring that even large datasets can be handled efficiently.

# 5. Data Sources

The data used in this project comes from **Myntra's fashion product database**, which includes a wide variety of fashion items, such as:

- **Product Descriptions**: A textual description of the product, which includes details such as fabric, style, usage, and more.
- **Metadata**: Additional data such as product ID, price, brand, and other relevant product attributes.

This data was pre-processed to convert product descriptions and features into text embeddings, which were then stored in **ChromaDB** for efficient search.

The dataset can be downloaded from the link

# 6. Key Design Decisions

## 1. Choice of Retrieval-Augmented Generation (RAG):

We chose the **RAG** approach because it allows the system to combine the best of **retrieval-based search** and **generation-based responses**. The retrieval phase ensures that the most relevant products are found, and the generation phase enriches these results with context and detail, making the system more interactive and informative.

## 2. Use of OpenAI Embeddings:

We opted for **OpenAI's embeddings** for their proven performance in capturing semantic meaning and their ability to generalize across various domains. This made them ideal for transforming the product descriptions and user queries into vector representations, facilitating accurate search results.

## 3. ChromaDB for Vector Storage:

**ChromaDB** was selected as the vector store due to its high performance and ability to scale effectively with large datasets. ChromaDB allows for fast retrieval times and supports similarity-based search efficiently, making it a good fit for this project.

## 4. LangChain for Orchestration:

We used **LangChain** as the framework to tie all components together. LangChain's flexibility and native support for LLMs and vector stores made it a natural choice for implementing the **RAG** workflow.

# 7. Challenges Encountered

### 1. Data Quality and Preprocessing:

One of the first challenges faced was ensuring that the product data from Myntra was clean and consistent. The data had to be pre-processed effectively to remove any inconsistencies and format it in a way that could be easily converted into embeddings. For instance, product descriptions had to be standardized to ensure consistent embedding generation.

### 2. Scaling the Vector Database:

As the dataset grew, ensuring that the **ChromaDB** could handle a large number of products without significant performance degradation became a challenge. This required careful optimization of the search index and parameters to balance between search speed and quality.

### 3. Fine-Tuning the Generation Model:

While the search phase returned relevant results, generating coherent and contextually rich descriptions was not always straightforward. It required some fine-tuning of the **OpenAI GPT model** to ensure that the generated text was helpful, accurate, and informative. Overfitting on specific product types was another issue that needed attention.

### 4. Handling Ambiguity in User Queries:

User queries can be ambiguous, and understanding the true intent behind vague queries like "stylish shoes" or "affordable jeans" was challenging. Implementing query refinement and disambiguation techniques, such as clarifying questions or expanding the query space, helped address this challenge.

# 8. Conclusion

The **Semantic Spotter Project** represents a significant leap forward in improving product search within the fashion domain. By combining **retrieval-based search** with **generation-based responses**, the system provides users with not only relevant products but also enriched, context-aware responses that enhance the shopping experience. Despite the challenges encountered—such as data preprocessing, vector store optimization, and query interpretation —the project successfully achieved its goals of improving search relevance, generating rich content, and creating a user-friendly experience.

This project lays the foundation for future enhancements, such as incorporating real-time data updates, refining the generation model for specific fashion categories, and enhancing the system's ability to handle more complex user queries. Ultimately, the **Semantic Spotter** can serve as a template for building advanced search systems in other domains as well.