



**ÇANKAYA UNIVERSITY
FACULTY OF ENGINEERING
COMPUTER ENGINEERING
DEPARTMENT**

**CENG 407
Project Report**

**Multimodal RAG-Based Product
Recommendation System**

Hazal KANTAR - 202111036
Ahmet Doğu kan GÜNDEMİR - 202111033
Ali Boran BEKTAS - 202111001
Hikmet Berkin BULUT - 202111057

Advisor: Assist. Prof. Dr. Serdar ARSLAN

Contents

Introduction	6
Work Plan	6
Literature Review	7
1 abstract	7
2 Introduction	7
3 What is a RAG System?	7
4 How Does Multimodal Work?	8
4.1 Image + Text Pairs	9
5 Multimodal RAG Architecture	9
6 Used Models	12
7 Deployment Challenges and Solutions	12
8 Top Papers in Multimodal Retrieval-Augmented Generation and Recommendation	13
9 Similar Systems	15
10 Conclusion	16
System Requirements Specification	17
11 Introduction	17
11.1 Purpose of This Document	17
11.2 Scope of This Project	17

12 General Description	17
12.1 Glossary (Definitions, Acronyms, and Abbreviations)	17
12.2 User Characteristics	18
12.3 Overview of Functional Requirements	18
12.4 General Constraints and Assumptions	19
13 Specific Requirements	20
13.1 Interface Requirements	20
13.1.1 User Interface	20
13.1.2 Hardware Interface	20
13.1.3 Software Interface	20
13.1.4 Communication Interfaces	21
13.2 Detailed Description of Functional Requirements	22
13.2.1 User Login and Registration	22
13.2.2 Personalized and Interactive Recommendation Generation	23
13.2.3 Trends Review	24
13.2.4 Real-Time Recommendation Updates	24
13.2.5 View Historical Recommendations	25
13.2.6 Profile Management	26
13.2.7 Search and Filter Options	27
13.3 Non-Functional Requirements	28
14 ANALYSIS - UML	30
14.1 Use Cases	30
14.1.1 Use Case Diagram	30
14.1.2 Description of Use Cases	31
14.2 Functional Modeling and Data Flow Diagrams (DFD)	38
14.2.1 Level - 0 Data Flow Diagram (Context Diagram)	38

14.2.2 Level - 1 Data Flow Diagram	39
15 CONCLUSION	39
System Design Document	40
16 Introduction	40
16.1 Purpose of this document	40
16.2 Definitions, Acronyms, and Abbreviations	40
17 System Overview	41
18 System Design	41
18.1 Architectural Design	41
18.1.1 Layered Architecture	42
18.1.2 Key Components	42
18.1.3 Benefits of Layered Architecture	42
18.2 Decomposition Description	43
18.2.1 Presentation Layer	43
18.2.2 Application Layer	44
18.2.3 Data Layer	45
18.2.4 External AI Services	46
18.2.5 Communication Between Layers	46
18.3 System Modeling	47
18.3.1 Activity Diagrams	47
18.3.2 Sequence Diagrams	50
18.3.3 Class Diagram	53
19 User Interface Design	54
19.1 User Interface of Profile Page	54

19.2 User Interface of Chatbot Screen	55
19.3 User Interface of Personal Recommendations Page	55
19.4 User Interface of Home Page	56
19.5 User Interface of Login Page	57
19.6 User Interface of Trending Page	58
Conclusion	58
References	60

Introduction

This project focuses on developing a Multimodal Retrieval-Augmented Generation (RAG) system to provide personalized product recommendations in the fashion and cosmetics industries. By using advanced AI technologies, the system combines user input, product images, and textual descriptions to generate accurate and relevant suggestions. Users can interact with the platform through a user-friendly web interface that supports both text and image-based chats.

The goal is to enhance the user experience by delivering recommendations that align with personal preferences, sustainability criteria, and current trends. The system integrates Large Language Models (LLMs) and external APIs, such as OpenAI, Gemini, and Hugging Face, to process multimodal data effectively. This approach allows for a more comprehensive understanding of user needs and market demands.

This report provides a comprehensive overview of the project, including the Literature Review, Project Work Plan, Software Requirements Specification (SRS), and Software Design Document (SDD), which detail the system's functionality, architecture, and design components. Together, these documents outline the project's objectives, technical specifications, and implementation roadmap.

Work Plan

The following image provides a detailed overview of the project work plan:

Multimodal-RAG-Based-Product-Recommendation-System		WEEKS														
Procedural Steps	Current State	1 30 September - 6 October	2 7 October - 13 October	3 14 October - 20 October	4 21 October - 27 October	5 28 October - 3 November	6 4 November - 10 November	7 11 November - 17 November	8 18 November - 24 November	9 25 November - 1 December	10 2 December - 8 December	11 9 December - 15 December	12 16 December - 22 December	13 23 December - 29 December	14 30 December - 5 January	15 6 January - 12 January
Team Setup	Completed															
Project Proposal Form	Completed															
Project Selection Form	Completed															
Github Repository	Completed															
Project Work Plan	Completed															
Literature Review	Completed															
Software Requirements Specification	Completed															
Project Webpage	Completed															
Software Design Description	Completed															
Project Report / Project Tracking Form	In Progress															
Presentation	In Progress															

Figure 1: Project Work Plan

Literature Review

1 abstract

In the rapidly evolving domains of fashion and cosmetics, personalized recommendations are crucial for enhancing user experience and driving sales. Traditional recommendation systems often rely solely on textual data, which may not fully capture fashion and cosmetic products' rich, multimodal nature. This literature review explores integrating Retrieval Augmented Generation (RAG) systems with multimodal data, specifically image and text pairs, to improve recommendation accuracy. We explore the architecture of multimodal RAG systems, discuss their application in product recommendations, examine the algorithms and models employed, and address deployment challenges. Additionally, we highlight top papers in the field and compare similar systems to comprehensively understand current advancements and future directions.

2 Introduction

The fashion and cosmetics industries are characterized by rapidly changing trends and a large number of products, making personalized recommendations essential for consumers navigating these markets. Traditional recommendation systems have primarily utilized textual data, such as product descriptions and user reviews, to infer user preferences. However, fashion and cosmetic products are inherently visual, and textual data alone may not be sufficient to capture their full essence. To address this limitation, integrating visual data with textual information has become a focal point in developing more effective recommendation systems. This integration is performed by Retrieval Augmented Generation systems, which combine retrieval mechanisms with generative models to enhance recommendation accuracy.

3 What is a RAG System?

Retrieval-augmented generation (RAG) systems are advanced models that enhance the capabilities of generative models by incorporating external knowledge through retrieval mechanisms. In traditional generative models, all knowledge is stored within the model's parameters, which can lead to limitations in scalability and adaptability. RAG systems address this by retrieving relevant information from external sources, such as databases or the internet, and integrating it into the generation process. This approach allows the model to produce more accurate and contextually relevant outputs, as it can access up-to-date information beyond its training data.

4 How Does Multimodal Work?

Multimodal machine learning integrates and processes information from multiple data modalities—such as text, images, audio, and video—to enhance a model’s understanding and performance. Each modality offers unique insights; for instance, text provides semantic context, while images deliver visual details. By combining these diverse data types, multimodal systems can capture a more comprehensive representation of the information, leading to improved accuracy and robustness in various tasks.

The process begins with the extraction of features from each modality using specialized encoders. For textual data, natural language processing techniques are employed, whereas convolutional neural networks are present abstractly in the back for image data. These encoders transform raw data into structured representations that encapsulate the essential characteristics of each modality. Once features are extracted, the challenge lies in effectively integrating them. This integration can occur at different stages:

- **Early Fusion:** Combines raw data from all modalities before feature extraction, allowing the model to learn joint representations from the outset.
- **Late Fusion:** Processes each modality independently through separate models and merges their outputs at a later stage, facilitating decision-level integration.
- **Hybrid Fusion:** Incorporates elements of both early and late fusion, enabling flexibility in how and when modalities are combined.

Advanced techniques, such as attention mechanisms, are often employed to dynamically weigh the importance of each modality’s contribution. This ensures that the model focuses on the most relevant information from each source, enhancing its decision-making capabilities.

In the context of Retrieval Augmented Generation systems, multimodal integration is particularly valuable. RAG systems retrieve relevant information from external sources to augment the generation process, improving the quality and relevance of the output. By incorporating multiple modalities, RAG systems can access a richer set of information, leading to more informed and contextually appropriate responses. For example, in fashion and cosmetics, a multimodal RAG system can analyze both textual descriptions and visual images of products to generate personalized recommendations or detailed product summaries. This comprehensive approach uses the strengths of each modality, resulting in a more nuanced and effective system.

In summary, multimodal machine learning works by extracting and integrating features from various data types, allowing models to leverage complementary information. This integration enhances the model’s ability to understand complex data, leading to improved performance across a wide range of applications.

4.1 Image + Text Pairs

In multimodal recommendation systems, image and text pairs are utilized to represent products comprehensively. The textual component may include product descriptions, specifications, and user reviews, while the visual component consists of product images. Combining these modalities allows the system to analyze both the semantic content of the text and the visual features of the images. This dual analysis enables the system to capture nuances such as style, color, and design, which are particularly important in fashion and cosmetics. For example, a user interested in a “red evening dress” would benefit from recommendations that consider both the textual description (“red,” “evening dress”) and the visual appearance of the dress.

5 Multimodal RAG Architecture

Multimodal Retrieval-Augmented Generation architectures enhance the capabilities of large language models by integrating information retrieval mechanisms that handle diverse data modalities, such as text and images. This integration enables models to generate more accurate and contextually relevant outputs by leveraging external knowledge sources. A typical Multimodal RAG system comprises two primary components:

Retriever: This component is responsible for fetching relevant information from external databases or knowledge bases. In a multimodal context, the retriever processes various data types, including text and images, to identify pertinent content. For instance, models like CLIP (Contrastive Language-Image Pre-training) are employed to generate embeddings for both textual and visual data, facilitating effective retrieval across modalities. [1]

Generator: After retrieval, the generator utilizes the fetched information to produce coherent and contextually enriched responses. Advanced language models such as Llama 3 and Gemini are often used in this role. These models are fine-tuned to handle multimodal inputs, enabling them to generate outputs that seamlessly integrate information from both text and images.

Query Processing: Upon receiving a user query, the system determines the nature of the input, which could be textual, visual, or a combination of both.

The interaction between these components operates as follows:

Multimodal Retrieval: The retriever processes the query to extract relevant information from the knowledge base. This involves generating embeddings for the query and comparing them with stored embeddings to identify the most relevant data. For example, CLIP can encode both the query and the knowledge base entries into a shared embedding space, allowing for efficient similarity matching.[1]

Contextual Generation: The generator then synthesizes the retrieved information to produce a response that is both accurate and contextually appropriate. Models like Llama 3 and Gemini are designed to handle such tasks, leveraging their advanced language

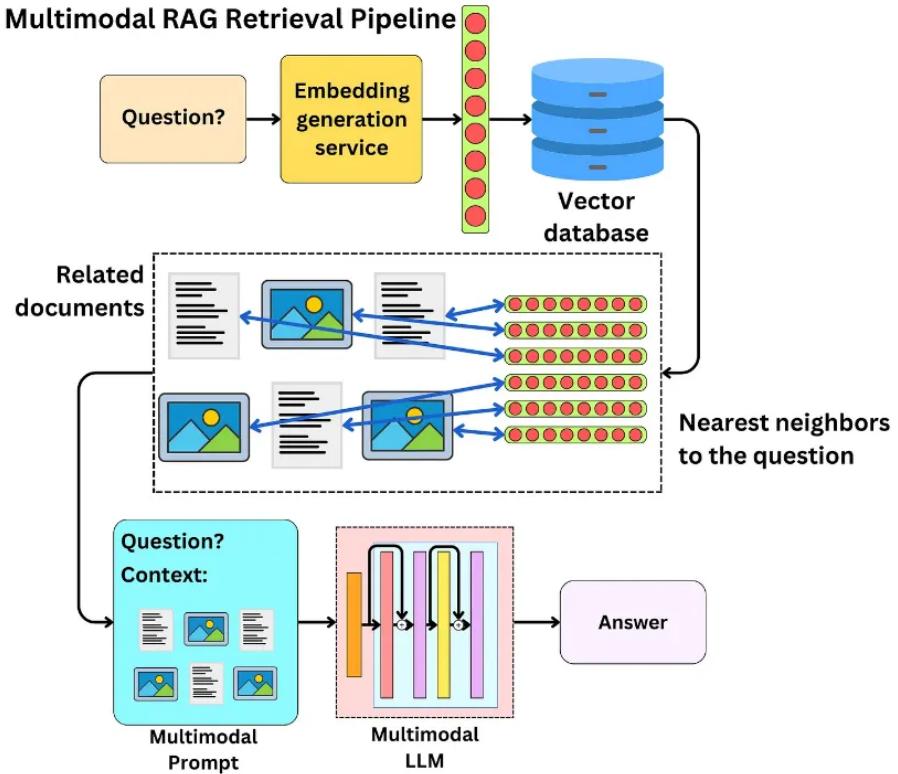


Figure 2: Retrieval Pipeline

understanding capabilities to generate coherent outputs.

This architecture is particularly beneficial in applications like image or product recommendation systems. By processing both textual descriptions and visual content, the system can provide more personalized and relevant recommendations. For instance, in the fashion and cosmetics industry, a Multimodal RAG system can analyze user preferences expressed through text and images to suggest products that align with current trends and individual tastes.

Implementing such systems involves several challenges, including the need for large-scale multimodal datasets, efficient retrieval mechanisms, and the integration of diverse models. However, advancements in models like CLIP, Llama 3, Gemini, and GPT-4 have significantly contributed to overcoming these challenges, leading to more sophisticated and effective Multimodal RAG systems.

In summary, Multimodal RAG architectures represent a significant advancement in AI, enabling systems to process and generate information across multiple modalities. By using models like CLIP for retrieval and Llama 3 or Gemini for generation, these systems can provide more accurate, contextually relevant, and personalized outputs, particularly in domains such as fashion and cosmetics.

Detailed Steps in Multimodal RAG:

- 1. User Query (Text + Image):** The process begins when a user submits a query using text, an image, or both. In a fashion use case, a user might submit an image of an

outfit with a query like, “Show me similar styles for this outfit.”

2. Text and Image Embeddings: Both the text and image are passed through encoders (e.g., CLIP). The text is processed by a Transformer-based text encoder, while the image goes through an image encoder (Vision Transformer). Both outputs are projected into the same shared embedding space.

3. Multimodal Fusion: The text and image embeddings are fused into a single query embedding, which combines information from both modalities. This fusion may be done via concatenation or attention-based mechanisms, allowing the system to fully capture the intent of the user’s multimodal query.

4. Retrieve Relevant Data (from Image and Text Databases): The fused query embedding is used to search a database that stores both text and image embeddings. The database retrieves images (e.g., visuals of outfits) and relevant text information, such as descriptions, metadata, or reviews.

5. Combine Retrieved Data (Text + Images): Once relevant data is retrieved, it combines both image and text data to generate a richer response. For example, a fashion recommendation system might return images of matching outfits along with style suggestions.

6. Generate Response (Text Output): The retrieved data is fed into a generative model like GPT-4 to create a conversational response. For instance, the system might return, “Here are accessories that would go well with your red dress: a pair of black heels and a silver necklace.”

Text/Image Embeddings - Encoders: Encoders are neural network models designed to transform input data into fixed-size vectors, known as embeddings, which capture the semantic meaning of the input.

- **Text Encoders:** These models process textual data to generate embeddings that represent the semantic content of the text. Transformer-based models, such as BERT or the text encoder component of CLIP, are commonly used. They analyze the input text, capturing contextual information and relationships between words to produce meaningful embeddings.
- **Image Encoders:** These models process visual data to generate embeddings that capture the semantic content of images. Convolutional Neural Networks (CNNs) or Vision Transformers (ViTs) are typically employed. They analyze the visual features of the image, such as shapes, colors, and textures, to produce embeddings that represent the image’s content.

In the context of Multimodal RAG, models like CLIP are utilized to generate embeddings for both text and images. CLIP employs a contrastive learning approach, training the model to align text and image embeddings in a shared space. This alignment enables the system to effectively retrieve and integrate information across modalities, enhancing the generation of contextually relevant responses.

6 Used Models

In the domain of fashion and cosmetics, integrating multimodal Retrieval-Augmented Generation systems has become pivotal for enhancing user experiences through personalized recommendations and trend analyses. These systems efficiently combine textual and visual data to generate contextually rich and accurate outputs. Several advanced models have been effective in this integration:

- **CLIP (Contrastive Language–Image Pretraining):** Developed by OpenAI, CLIP is a foundational model that aligns textual and visual representations by training on a vast dataset of image-text pairs. This alignment enables the model to understand and generate content that seamlessly integrates both modalities, making it particularly effective for applications requiring a deep understanding of visual and textual data.
- **LLaMA 3 (Large Language Model Meta AI):** LLaMA 3 is a state-of-the-art language model designed to process and generate human-like text. Its advanced architecture allows it to comprehend complex textual inputs and produce coherent and contextually relevant outputs, which is essential for generating detailed product descriptions and fashion trend analyses.
- **Gemini:** Gemini is a multimodal model that integrates both language and vision capabilities. It processes and generates content that encompasses textual and visual elements, making it suitable for applications like image-based product recommendations and visual content generation.
- **GPT-4:** As one of OpenAI's most advanced language models, GPT-4 excels in understanding and generating human-like text. Its capabilities are enhanced when integrated with multimodal systems, allowing it to provide detailed and contextually rich responses that incorporate both textual and visual information.

By utilizing these models, multimodal RAG systems in the fashion and cosmetics industry can deliver more personalized and accurate recommendations, effectively combining the strengths of both textual and visual data processing.

7 Deployment Challenges and Solutions

Deploying Multimodal Retrieval-Augmented Generation systems, especially in domains like fashion and cosmetics, presents several challenges. These challenges include handling diverse data modalities, ensuring system scalability, maintaining data privacy, and managing computational resources. Addressing these issues is crucial for the effective implementation of multimodal RAG systems.

1. Handling Diverse Data Modalities

Challenge: The fashion and cosmetics industries rely heavily on both visual and textual data. Integrating these modalities into a cohesive system is complex.

Solution: Utilize models like CLIP (Contrastive Language–Image Pre-training) to create unified embeddings for images and text, facilitating seamless integration. Additionally, employing multimodal models such as GPT-4V and Gemini can enhance the system's ability to process and generate content across different data types.

2. Ensuring System Scalability

Challenge: As the volume of data grows, the system must efficiently handle increased loads without compromising performance.

Solution: Implement scalable architectures using cloud services like Azure or AWS, which offer flexible resources to accommodate varying workloads. Incorporating vector databases, such as Milvus, or ChromaDB, can optimize the storage and retrieval of embeddings, ensuring quick access to relevant information.

3. Maintaining Data Privacy

Challenge: Handling sensitive customer data necessitates strict adherence to privacy regulations.

Solution: Implement robust data anonymization techniques and comply with data protection laws. Regular checks and the use of secure data storage solutions are essential to safeguard user information.

4. Ensuring Real-time Performance

Challenge: Users expect quick responses, making latency a critical factor.

Solution: Optimize retrieval algorithms and implement efficient caching mechanisms to reduce response times. Employing asynchronous processing can further enhance system responsiveness.

8 Top Papers in Multimodal Retrieval-Augmented Generation and Recommendation

The integration of multimodal data into retrieval-augmented generation systems has been a focus of recent research. Below are papers that have significantly advanced this field:

1. “MuRAG: Multimodal Retrieval-Augmented Generator for Open Question Answering over Images and Text”

Authors: Wenhui Chen, Hexiang Hu, Xi Chen, Pat Verga, William W. Cohen

Summary: This paper introduces MuRAG, a model that enhances language generation by accessing an external multimodal memory comprising images and text. Pre-trained on large-scale image-text and text-only datasets, MuRAG demonstrates state-of-the-art performance in open-domain question-answering tasks requiring multimodal reasoning. [2]

2. “Retrieving Multimodal Information for Augmented Generation: A Survey”

Authors: Ruochen Zhao, Hailin Chen, Weishi Wang, et al.

Summary: This comprehensive survey reviews methods that assist generative models by retrieving multimodal knowledge, including images, code, tables, graphs, and audio. It discusses the applications, challenges, and future directions of multimodal retrieval-augmented generation. [3]

3. “Retrieval-Augmented Multimodal Language Modeling”

Authors: Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, et al.

Summary: The authors propose a retrieval-augmented multimodal model that enables a base multimodal model to refer to relevant text and images fetched from external memory. This approach enhances text-to-image and image-to-text generation tasks. [4]

4. “MLLM Is a Strong Reranker: Advancing Multimodal Retrieval-Augmented Generation via Knowledge-Enhanced Reranking and Noise-Injected Training”

Authors: Zhanpeng Chen, Chengjin Xu, Yiyan Qi, Jian Guo

Summary: This paper presents a framework that improves multimodal retrieval-augmented generation by incorporating knowledge-enhanced reranking and noise-injected training, leading to more robust and accurate multimodal generation. [5]

5. “I Want This Product but Different: Multimodal Retrieval with Synthetic Query Expansion”

Authors: Ivona Tautkute, Tomasz Trzcinski

Summary: The authors address the problem of media retrieval using multimodal queries by proposing a framework that expands the query with a synthetically generated image, capturing semantic information from both image and text inputs. [6]

6. “MuRAR: A Simple and Effective Multimodal Retrieval and Answer Refinement Framework for Multimodal Question Answering”

Authors: Zhengyuan Zhu, Daniel Lee, Hong Zhang, et al.

Summary: MuRAR enhances text-based answers by retrieving relevant multimodal data and refining responses to create coherent multimodal answers, demonstrating improved performance in multimodal question-answering tasks. [7]

7. “Retrieving Multimodal Information for Augmented Generation: A Survey”

Authors: Ruochen Zhao, Hailin Chen, Weishi Wang, et al.

Summary: This survey provides an in-depth review of methods that assist generative models by retrieving multimodal knowledge, and discussing their applications and future directions. [8]

9 Similar Systems

In the domains of fashion and cosmetics, several systems have been developed to provide personalized recommendations by integrating user preferences and multimodal data. However, most existing solutions focus on either fashion or cosmetics individually, rather than combining both to deliver a unified, user-specific experience.

1. Fashion Personalization Systems:

- Intelistyle: This platform offers personalized fashion recommendations by analyzing user preferences and current trends. It utilizes AI to suggest outfits and styles tailored to individual tastes. [9]
- The Iconic: An online retailer that combines AI technology with fashion and beauty, offering personalized recommendations to enhance user experience [10]

2. Cosmetic Personalization Systems:

- Perfect Corp: Known for its AI and AR solutions, Perfect Corp collaborates with beauty brands to provide personalized product recommendations and virtual try-on experiences. [11]
- Formulate: A personalized haircare brand that creates customized products based on individual hair needs and preferences. [12]

3. Integrated Fashion and Cosmetic Systems:

- While some platforms offer both fashion and cosmetic products, they often treat these categories separately without a cohesive integration. For instance, Cover-Girl's flagship store provides personalized makeup try-on experiences but does not integrate fashion recommendations. [13]

4. Our Approach: Our system differentiates itself by seamlessly combining fashion and cosmetic recommendations into a single, cohesive output. By analyzing user inputs and preferences, we provide tailored suggestions that encompass both clothing and beauty products, ensuring a harmonious and personalized ensemble. This integrated approach enhances user satisfaction by delivering comprehensive recommendations that align with individual styles and preferences.

10 Conclusion

The integration of multimodal retrieval-augmented generation (RAG) systems within the fashion and cosmetics industries signifies a transformative advancement in personalized user experiences. By using diverse data modalities such as text, images, and user preferences these systems offer tailored recommendations that align closely with individual tastes and needs.

The development of sophisticated architectures, including models like CLIP, LLaMA 3, and Gemini, has been pivotal in enhancing the capabilities of multimodal RAG systems. These models facilitate the seamless fusion of various data types, enabling more accurate and contextually relevant outputs. For instance, CLIP’s ability to understand and generate both text and images allows for more nuanced product recommendations, while LLaMA 3’s advanced language modeling contributes to more coherent and personalized interactions.

Despite these advancements, deploying multimodal RAG systems presents challenges, such as managing computational resources, ensuring data privacy, and maintaining system scalability. Addressing these issues requires a combination of robust infrastructure, efficient algorithms, and adherence to ethical standards.

The fusion of fashion and cosmetics recommendations within a single output, tailored to user-specific inputs, represents a significant innovation. This approach not only enhances user satisfaction by providing comprehensive style suggestions but also sets a new standard for personalized services in the industry.

In summary, the evolution of multimodal RAG systems, supported by advanced models and thoughtful integration strategies, is reshaping the landscape of fashion and cosmetics recommendations. As these technologies continue to grow, they deliver increasingly personalized and engaging experiences to users, thereby driving innovation in the industry.

System Requirements Specification

11 Introduction

11.1 Purpose of This Document

The purpose of this document is to outline the software requirements for the development of a Multimodal Retrieval-Augmented Generation (RAG)-Based Product Recommendation System. This document serves as a reference for stakeholders, including project managers, developers, designers, and testers, ensuring a shared understanding of the project's goals, features, constraints, and technical requirements. It provides a comprehensive framework for system implementation, verification, and maintenance.

11.2 Scope of This Project

This project focuses on designing an intelligent recommendation platform for the fashion and cosmetics domains, using multimodal data sources such as textual descriptions, product images, sustainability certifications, and trend analytics. By utilizing advanced AI technologies, including LLMs and multimodal embeddings, the system provides eco-conscious and trend-aligned product recommendations tailored to individual user preferences. The platform consists of a Flask-based backend, a React-powered frontend, and a vector database for efficient data handling. Its ultimate goal is to promote sustainable consumer choices through a user-friendly, adaptive, and visually engaging interface.

12 General Description

12.1 Glossary (Definitions, Acronyms, and Abbreviations)

- **AI:** Artificial Intelligence.
- **API:** Application Programming Interface, a connection between computers or between computer programs.
- **CLIP:** Contrastive Language–Image Pretraining, a model that creates unified embeddings for text and images.[16]
- **Flask:** A lightweight Python web framework for backend API development.[17]
- **Frontend:** The part of the application that interacts directly with the user.
- **Backend:** The part of an application that is not directly accessed by the user, typically responsible for storing and manipulating data.
- **Gemini:** A suite of AI models developed by Google.[18]

- **Hugging Face:** A company providing natural language processing tools and models.[19]
- **HCI:** Human-Computer Interaction; the study of how people interact with computers and to design technologies that let humans interact with computers in novel ways.
- **LLM:** Large Language Model, a type of computational model designed for natural language processing tasks.
- **LLaMA:** Large Language Model Meta AI; a family of large language models developed by Meta AI.[20]
- **OpenAI:** An AI research and deployment company.
- **RAG:** Retrieval-Augmented Generation, A framework combining retrieval-based methods with generative models for improved contextual output.
- **React:** A JavaScript library for building user interfaces.
- **Vector Database:** A database optimized for storing and retrieving vector embeddings.

12.2 User Characteristics

The primary users of this system are environmentally conscious consumers seeking personalized recommendations for fashion and cosmetic products. They are comfortable with digital platforms and expect seamless, intuitive interactions. Secondary users include fashion retailers and brands interested in promoting sustainable products to a targeted audience.

12.3 Overview of Functional Requirements

The system involves a variety of functionalities to provide a seamless user experience and effective personalized recommendations. It includes user login and registration, enabling secure access to the platform while storing individual preferences and interaction history. Personalized and interactive recommendation generation lies at the core of the system, using multimodal data and advanced AI models to create suggestions tailored to the user's preferences. Additionally, users can review ongoing trends in the fashion and cosmetics domains, offering insights into popular and eco-conscious products.

Real-time recommendation updates ensure that users receive the most relevant suggestions based on current trends and new data. Historical recommendations are accessible through the platform, allowing users to revisit previous suggestions for convenience. Profile management features allow users to update their preferences, adjust their interaction history, and manage saved recommendations. Furthermore, a robust search and filtering mechanism is integrated to help users explore the product catalog more efficiently, ensuring a user-centric and dynamic recommendation platform.

12.4 General Constraints and Assumptions

Constraints

- **Data Availability:** The system's performance relies on continuous access to current and accurate product data, including images, descriptions, availability, and sustainability certifications.
- **Performance Limitations:** The efficiency and responsiveness of the recommendation engine are influenced by the capabilities of the employed AI models and the computational resources at hand, which may affect response times and recommendation quality.
- **Regulatory Compliance:** The system must adhere to data privacy regulations and ethical guidelines, ensuring responsible and transparent handling of user data.
- **Scalability:** Accommodating a growing user base and expanding product catalogs requires efficient data processing and storage solutions to maintain system performance.
- **Compatibility:** The system must ensure compatibility with a range of web browsers, including Chrome, Firefox, Safari, and Edge, to support a diverse user base.
- **Accessibility:** As an online platform, the system must provide uninterrupted access and maintain optimal functionality for users across different devices and network conditions.

Assumptions

- **User Connectivity:** It is assumed that users have reliable internet access and utilize modern web browsers to interact with the platform.
- **Accurate User Input:** The effectiveness of personalized recommendations depends on users providing precise and comprehensive preference information.
- **Retailer Collaboration:** Successful operation relies on effective collaboration with retailers to ensure the availability and accuracy of product sustainability information, which is critical for delivering environmentally conscious recommendations.
- **User Proficiency:** Users are presumed to possess a basic level of digital literacy, enabling them to navigate the platform and utilize its features effectively.
- **Stable Operating Environment:** The system is expected to operate within a stable technological environment, with minimal disruptions due to software updates or hardware failures.

13 Specific Requirements

13.1 Interface Requirements

13.1.1 User Interface

The system shall provide an intuitive and responsive user interface accessible via modern web browsers. The frontend will be developed using React to ensure dynamic content rendering and a seamless user experience. Users will interact with the platform through various components, including:

- **Dashboard:** Displays personalized recommendations, recent trends, and user activity summaries.
- **Search and Filter Panel:** Allows users to search for products and apply filters based on categories, sustainability certifications, price ranges, and popularity.
- **Product Detail View:** Provides comprehensive information about selected products, including images, descriptions, sustainability credentials, and user reviews.
- **User Profile Management:** Enables users to view and edit personal information, manage preferences, and review historical recommendations.

The interface will adhere to accessibility standards to accommodate users with varying needs, ensuring a user-friendly experience for all.

13.1.2 Hardware Interface

The system will integrate with various software components and services to deliver its functionalities:

- **User Devices:** Must support modern web browsers compatible with the React-based frontend.
- **Server Infrastructure:** The backend, developed using Flask, will be hosted on cloud servers (e.g., AWS, Google Cloud Platform, Microsoft Azure) to ensure scalability and reliability.

13.1.3 Software Interface

The system will integrate with various software components and services to deliver its functionalities:

- **Backend Services:** The Flask-based backend will handle API requests, manage business logic, and interface with the database.

- **Database:** A vector database (e.g., ChromaDB, Milvus) will store multimodal embeddings for efficient retrieval.
- **External APIs:** Integration with external services such as OpenAI API, Gemini, and Hugging Face for embedding services will enhance recommendation accuracy.
- **Authentication Services:** Implementation of secure authentication protocols to manage user login and registration.

These software interfaces will be designed to ensure seamless communication between components, maintaining data integrity and system performance.

13.1.4 Communication Interfaces

Client-Server Communication The frontend (React-based) will communicate with the backend (Flask-based) via RESTful APIs. These APIs will manage data retrieval, user requests, and recommendations. All communication will be secured using HTTPS to protect data integrity and confidentiality.

Backend-Database Interaction The backend will interact with the vector database (e.g., ChromaDB, Milvus) using database-specific APIs and libraries. These interactions will include storing, updating, and retrieving multimodal embeddings for recommendation generation.

Integration with External APIs The system will communicate with third-party services like OpenAI API, Gemini, and Hugging Face to access advanced AI models and embedding services. Calls to these APIs will be managed via asynchronous methods to optimize performance and ensure responsiveness.

User Notifications The system will include communication interfaces for delivering notifications to users, such as alerts for new trends, recommendations, or updates. Notifications will be sent through web-based mechanisms such as in-app messages or push notifications.

Error Handling and Logging Communication interfaces will include robust error-handling mechanisms to manage issues such as failed API calls, network disruptions, or database inconsistencies. A logging system will capture detailed logs of communication activities to support debugging and system maintenance.

13.2 Detailed Description of Functional Requirements

13.2.1 User Login and Registration

Name	User Login and Registration
Purpose/Description	Enables users to create accounts, securely log in, and access personalized features of the platform.
Inputs	<ul style="list-style-type: none">• User-provided email, username, and password during registration.• Login credentials for authentication.
Processing	<ul style="list-style-type: none">• Validation of input fields to ensure compliance with security policies (e.g., strong passwords).• Storage of user credentials using encrypted formats.• Authentication via secure protocols.
Outputs	<ul style="list-style-type: none">• Confirmation of successful registration or login.• Error messages for failed attempts (e.g., incorrect credentials or duplicate accounts).
Error Handling	<ul style="list-style-type: none">• Display specific errors for invalid inputs or registration conflicts.• Account recovery options in case of forgotten passwords.

13.2.2 Personalized and Interactive Recommendation Generation

Name	Personalized and Interactive Recommendation Generation
Purpose/Description	Generates personalized product recommendations based on user preferences, interactions, and historical data.
Inputs	<ul style="list-style-type: none"> User preferences, profile data, and interaction history. Multimodal product data (e.g., text, images, sustainability metrics).
Processing	<ul style="list-style-type: none"> Retrieval of relevant data from the vector database using user embeddings. Integration of results from AI models (e.g., RAG pipeline) to provide personalized recommendations. Interactive updates based on user feedback (e.g., refining recommendations).
Outputs	<ul style="list-style-type: none"> A list of products tailored to user preferences. Real-time updates to recommendations when user preferences are adjusted.
Error Handling	<ul style="list-style-type: none"> Data Gaps: Provide generic recommendations when user data is incomplete. User Feedback: Allow users to report inaccurate or irrelevant recommendations. Fallback recommendations in case of model or API failure.

13.2.3 Trends Review

Name	Trends Review
Purpose/Description	Allows users to review current trends in fashion and cosmetics, including eco-friendly options and popular products.
Inputs	<ul style="list-style-type: none"> External Data Sources: Trend data from APIs like NewsAPI, social media platforms, and industry reports. Internal Product Data: Details of products related to trending keywords, categories, or styles. User Preferences: Personal interests or saved categories influencing trend presentation.
Processing	<ul style="list-style-type: none"> Aggregation and filtering of trend data based on user preferences and sustainability criteria. Visualization of trend summaries (e.g., trending product categories, seasonal highlights).
Outputs	<ul style="list-style-type: none"> Trend Highlights: A list of popular trends, including associated products and descriptions. Visual Summaries: Graphs showing the evolution and impact of trends over time.
Error Handling	<ul style="list-style-type: none"> Data Source Failures: Display a fallback message with alternative suggestions. Analysis Errors: Log issues in trend extraction and ensure retrying data processing.

13.2.4 Real-Time Recommendation Updates

Name	Real-Time Recommendation Updates
Purpose/Description	Provides users with up-to-date product recommendations that reflect the latest trends, availability, and user interactions, ensuring relevance and timeliness in suggestions.

Inputs	<ul style="list-style-type: none"> User Interactions: Real-time data on user behaviors, such as clicks, views, purchases, and feedback. Product Data: Continuous updates on product details, including availability, pricing, and sustainability certifications. Trend Information: Current data on fashion and cosmetic trends sourced from industry reports, social media, and market analyses.
Processing	<ul style="list-style-type: none"> Data Aggregation: Collect and integrate real-time inputs from various sources to form a comprehensive dataset. Dynamic Analysis: Utilize AI algorithms to analyze aggregated data, identifying patterns and shifts in user preferences and market trends. Recommendation Adjustment: Modify existing recommendations based on new insights, ensuring alignment with the most recent data.
Outputs	<ul style="list-style-type: none"> Updated Recommendations: Present users with a refreshed list of product suggestions that mirror current trends and personal preferences. User Notifications: Alert users to significant updates or changes in recommendations.
Error Handling	<ul style="list-style-type: none"> Data Latency Management: Handle delays in data updates to ensure smooth functioning without noticeable lags. Fallback Strategies: Revert to the most recent stable data if real-time data is unavailable or delayed. User Feedback Integration: Monitor and address user-reported issues for continuous refinement.

13.2.5 View Historical Recommendations

Name	View Historical Recommendations
Purpose/Description	Enables users to view previously recommended products and revisit their preferences.
Inputs	<ul style="list-style-type: none"> User Profile Data: Historical records of recommendations generated for the user, stored in the system database.
Processing	<ul style="list-style-type: none"> Data Retrieval: Query the database to fetch historical recommendation records linked to the user profile. Organization and Sorting: Arrange the retrieved records based on relevance, time, or user-defined criteria for easy navigation.
Outputs	<ul style="list-style-type: none"> Historical Recommendations Display: A user-friendly visualization of past recommendations, categorized and sorted for clarity.
Error Handling	<ul style="list-style-type: none"> Data Unavailability: Notify the user if no historical data is available and provide guidance on generating new recommendations. Database Connectivity Issues: Display an error message and retry option if there is a temporary failure in fetching data. Search and Filter Errors: Offer default views if user-defined filters return no results.

13.2.6 Profile Management

Name	Profile Management
Purpose/Description	Allows users to manage their personal details, preferences, and account settings.
Inputs	<ul style="list-style-type: none"> User-Provided Data: Updates to personal information, preferences, and notification settings.

Processing	<ul style="list-style-type: none"> • Data Validation: Verify the accuracy and completeness of the new information provided by the user. • Profile Updates: Modify and save changes to user data in the database, ensuring secure storage and real-time synchronization. • Preference Integration: Adjust system settings and recommendation algorithms based on updated user preferences.
Outputs	<ul style="list-style-type: none"> • Updated Profile Information: Confirmation of changes made to user details and settings. • Enhanced Recommendations: Reflect updated preferences in the personalized recommendations provided by the system. • Notification Settings: Configure and display preferences for receiving updates or alerts.
Error Handling	<ul style="list-style-type: none"> • Input Errors: Prompt users to correct incomplete or invalid data entries. • Database Issues: Display a message if profile updates cannot be saved temporarily, and retry once the issue is resolved.

13.2.7 Search and Filter Options

Name	Search and Filter Options
Purpose/Description	Ensures users can locate items of interest by applying criteria such as categories, brands, and sustainability certification.

Inputs	<ul style="list-style-type: none"> Search Query: Keywords entered by the user to locate specific products. Filter Criteria: User-defined parameters, such as price range, brand, sustainability certifications, and trending status.
Processing	<ul style="list-style-type: none"> Query Matching: Search the database for products matching the user's input keywords and filter criteria. Filter Application: Refine search results based on selected attributes, such as relevance, price, or category.
Outputs	<ul style="list-style-type: none"> Search Results: Display a list of products that match the user's search query and filter preferences. Filter Summary: Show the applied filters alongside the results, allowing users to adjust them for refined searches.
Error Handling	<ul style="list-style-type: none"> No Results Found: Notify the user if no products match the search query or filters. Invalid Inputs: Prompt users to correct incomplete or inappropriate search queries or filter selections. System Delays: Display a loading indicator if the search process takes longer than expected.

13.3 Non-Functional Requirements

Backend Server Constraints

The backend server must efficiently handle multiple concurrent requests to ensure smooth functionality during high-traffic periods. The database queries shall be optimized using indexing and caching techniques, minimizing data retrieval time to under 50 milliseconds per query on average. This ensures seamless interaction with real-time recommendation generation, even under significant user loads.

System Responsiveness

The system shall generate product recommendations for user queries within an average response time of 5 seconds. This responsiveness will be maintained by employing asynchronous processing and optimized API calls, ensuring user satisfaction and minimizing wait times.

Availability

The system shall maintain at least 95 percent uptime, ensuring consistent availability for users. This high availability will be achieved through the use of cloud-based infrastructure with failover mechanisms and regular monitoring of system health to preemptively address potential issues.

Ease of Use

The user interface shall be intuitive and user-friendly, requiring no more than 10 minutes for a new user to understand the main features. This will be achieved through clean design principles, tooltips, and a comprehensive help section integrated into the web platform.

Privacy

The system shall comply with GDPR standards for handling user data. It will ensure the confidentiality of personal information and allow users to access, modify, or delete their data upon request. Security measures, such as data encryption and secure login protocols, will be implemented to safeguard user information.

Compatibility

The system shall support the latest versions of major web browsers, including Chrome, Firefox, Safari, and Edge, ensuring broad accessibility. Cross-platform compatibility tests will be conducted to identify and address any discrepancies, guaranteeing a consistent user experience.

Ethical Compliance and Bias Mitigation

The recommendation algorithms shall actively monitor and minimize biases in product suggestions, especially concerning gender, ethnicity, or body type. Regular audits of the recommendation engine will be conducted to detect and rectify any unintentional biases, fostering a fair and inclusive platform.

Transparency

The system shall provide clear explanations for its recommendations, enabling users to understand the rationale behind suggested products. A “Why this recommendation?” feature will be integrated into the user interface, offering insights into how user preferences, trends, and sustainability metrics influence the suggestions.

14 ANALYSIS - UML

14.1 Use Cases

14.1.1 Use Case Diagram

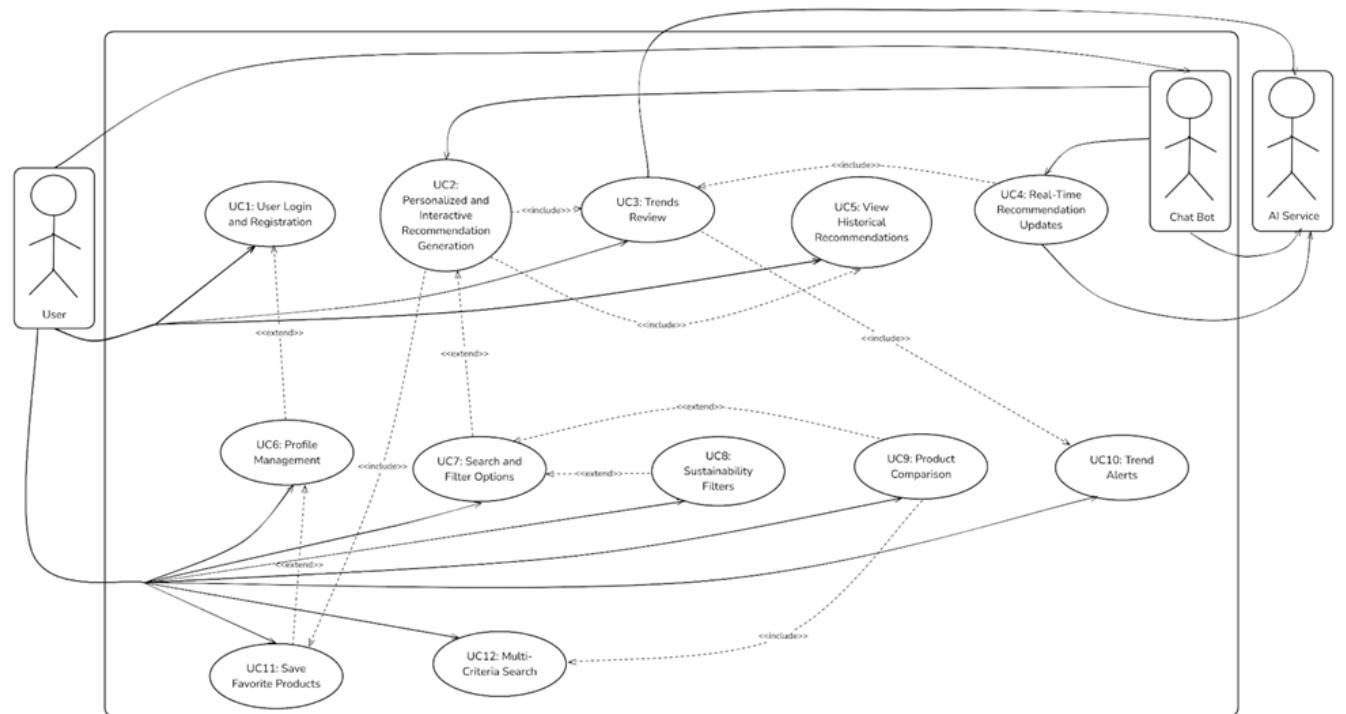


Figure 3: Main Use Case Diagram

14.1.2 Description of Use Cases

Use Case 1: User Login and Registration

Use Case Name	User Login and Registration
Use Case Number	1
Actors	User, System
Description	Users can create an account or log in to access personalized recommendations and manage their profile.
Precondition	The user accesses the system interface.
Scenario	<ol style="list-style-type: none"> 1. The user navigates to the login/registration page. 2. The user enters their credentials or registers by providing required details. 3. The system authenticates the user and grants access.
Postcondition	The user is successfully logged into the system or registered as a new user.
Exceptions	<ol style="list-style-type: none"> 1. Invalid login credentials. 2. User already exists during registration. 3. Network issues preventing communication with the backend.
Related Use Cases	Profile Management (UC6)

Use Case 2: Personalized and Interactive Recommendation Generation

Use Case Name	Personalized and Interactive Recommendation Generation
Use Case Number	2
Actors	User, System
Description	The system generates personalized recommendations based on user preferences, history, and current trends.
Precondition	The user has completed their profile and logged into the system.

Scenario	<ol style="list-style-type: none"> 1. The user logs into their account. 2. The system analyzes user preferences and history. 3. Recommendations are displayed based on real-time trends.
Postcondition	The user receives personalized recommendations tailored to their preferences and trends.
Exceptions	<ol style="list-style-type: none"> 1. Insufficient data in the user profile to generate meaningful recommendations. 2. System errors in retrieving or processing trend data.
Related Use Cases	Search and Filter Options (UC7), Save Favorite Products (UC11)

Use Case 3: Trends Review

Use Case Name	Trends Review
Use Case Number	3
Actors	User, System
Description	Users can view popular trends in fashion and cosmetics based on aggregated data from news sources, social media, and product sales.
Precondition	The system has up-to-date trend data available.
Scenario	<ol style="list-style-type: none"> 1. The user accesses the trends page. 2. The system displays real-time trends. 3. The user interacts with the data for insights.
Postcondition	The user gains insights into current fashion and cosmetics trends.
Exceptions	<ol style="list-style-type: none"> 1. Trend data is outdated or unavailable. 2. The system encounters issues while aggregating data from external APIs.
Related Use Cases	Personalized and Interactive Recommendation Generation (UC2), Real-Time Recommendation Updates (UC4)

Use Case 4: Real-Time Recommendation Updates

Use Case Name	Real-Time Recommendation Updates
Use Case Number	4
Actors	User, System
Description	The system updates recommendations dynamically as user preferences or external trend data changes.
Precondition	User preferences and trend data are updated in the system.
Scenario	<ol style="list-style-type: none"> 1. The system monitors changes in trends or user interactions. 2. Recommendations are updated dynamically. 3. Updated recommendations are displayed to the user.
Postcondition	The user receives real-time recommendations that reflect the latest preferences and trends.
Exceptions	<ol style="list-style-type: none"> 1. Delay in receiving updated trend data. 2. System overload due to high request volume.
Related Use Cases	Trends Review (UC3), Personalized and Interactive Recommendation Generation (UC2)

Use Case 5: View Historical Recommendations

Use Case Name	View Historical Recommendations
Use Case Number	5
Actors	User, System
Description	Users can view their past recommendations to revisit previous suggestions or track changes in trends over time.
Precondition	The system has a log of past recommendations.
Scenario	<ol style="list-style-type: none"> 1. The user navigates to their history page. 2. The system retrieves and displays historical recommendations.
Postcondition	The user views a history of their recommendations.

Exceptions	<ol style="list-style-type: none"> 1. Insufficient data in the history log. 2. System errors in retrieving historical data.
Related Use Cases	Personalized and Interactive Recommendation Generation (UC2)

Use Case 6: Profile Management

Use Case Name	Profile Management
Use Case Number	6
Actors	User, System
Description	Users can manage their personal details, preferences, and settings to customize their experience.
Precondition	The user is logged into their account.
Scenario	<ol style="list-style-type: none"> 1. The user navigates to the profile section. 2. The user modifies preferences or updates personal details. 3. The system saves the updated preferences.
Postcondition	User preferences and settings are successfully updated in the system.
Exceptions	<ol style="list-style-type: none"> 1. System fails to save the updated data due to database errors. 2. User inputs invalid data formats (e.g., non-numeric phone numbers).
Related Use Cases	User Login and Registration (UC1)

Use Case 7: Search and Filter Options

Use Case Name	Search and Filter Options
Use Case Number	7
Actors	User, System
Description	Users can search and filter products using criteria like price, color, brand, and sustainability attributes.
Precondition	The system has products and metadata for filtering.

Scenario	<ol style="list-style-type: none"> 1. The user enters search criteria in the filter options. 2. The system retrieves products matching the criteria. 3. Results are displayed in real-time.
Postcondition	The user views products filtered by their selected criteria.
Exceptions	<ol style="list-style-type: none"> 1. No products match the selected criteria. 2. Errors in retrieving product data due to system or database issues.
Related Use Cases	Personalized and Interactive Recommendation Generation (UC2), Product Comparison (UC9), Sustainability Filters (UC8)

Use Case 8: Sustainability Filters

Use Case Name	Sustainability Filters
Use Case Number	8
Actors	User, System
Description	Users can apply sustainability filters (e.g., eco-friendly materials, cruelty-free certifications) to refine product recommendations.
Precondition	The system contains sustainability data for the available products.
Scenario	<ol style="list-style-type: none"> 1. The user accesses the filter options. 2. The user selects one or more sustainability criteria. 3. The system retrieves and displays filtered results.
Postcondition	Only products meeting the selected sustainability criteria are displayed to the user.
Exceptions	<ol style="list-style-type: none"> 1. No products meet the selected criteria. 2. Incorrect filtering due to incomplete data in the database.
Related Use Cases	Search and Filter Options (UC7), Personalized and Interactive Recommendation Generation (UC2)

Use Case 9: Product Comparison

Use Case Name	Product Comparison
Use Case Number	9
Actors	User, System
Description	Users can compare multiple products based on attributes such as price, popularity, and sustainability.
Precondition	The user has added products to the comparison list.
Scenario	<ol style="list-style-type: none"> 1. The user selects products for comparison. 2. The system retrieves the attributes of the selected products. 3. The system displays a comparison table.
Postcondition	The user views a detailed comparison of the selected products.
Exceptions	<ol style="list-style-type: none"> 1. Insufficient product data for comparison. 2. System errors while retrieving product attributes.
Related Use Cases	Multi-Criteria Search (UC12), Search and Filter Options (UC7)

Use Case 10: Trend Alerts

Use Case Name	Trend Alerts
Use Case Number	10
Actors	User, System
Description	Users can subscribe to alerts for new trends in the fashion and cosmetics industry based on their preferences.
Precondition	The user has an active subscription and trend data is available.
Scenario	<ol style="list-style-type: none"> 1. The user enables trend alerts in their profile settings. 2. The system monitors trend data. 3. Alerts are sent via notifications or emails.
Postcondition	The user is notified about relevant trends.

Exceptions	<ol style="list-style-type: none"> 1. Trend data unavailable due to API errors. 2. Notifications fail to send due to connectivity issues.
Related Use Cases	Trends Review (UC3)

Use Case 11: Save Favorite Products

Use Case Name	Save Favorite Products
Use Case Number	11
Actors	User, System
Description	Users can save products to their favorites list for easy access and future reference.
Precondition	The user is logged into their account.
Scenario	<ol style="list-style-type: none"> 1. The user clicks the favorite icon for a product. 2. The system saves the product to the user's favorites list. 3. The user accesses their favorites later.
Postcondition	The selected product is added to the user's favorites list.
Exceptions	<ol style="list-style-type: none"> 1. Errors saving the product to the favorites list due to database issues. 2. The system fails to retrieve the favorites list when requested.
Related Use Cases	Personalized and Interactive Recommendation Generation (UC2), Profile Management (UC6)

Use Case 12: Multi-Criteria Search

Use Case Name	Multi-Criteria Search
Use Case Number	12
Actors	User, System
Description	Users can perform searches by combining multiple criteria, such as price range, brand, product category, and sustainability certifications.
Precondition	The system has a comprehensive database of product metadata and attributes.

Scenario	<ol style="list-style-type: none"> 1. The user inputs multiple criteria in the search form. 2. The system processes the criteria and retrieves matching products. 3. Results are displayed.
Postcondition	The user views a filtered list of products matching the selected criteria.
Exceptions	<ol style="list-style-type: none"> 1. No products match all the selected criteria. 2. Errors in data retrieval or processing.
Related Use Cases	Search and Filter Options (UC7), Product Comparison (UC9)

14.2 Functional Modeling and Data Flow Diagrams (DFD)

14.2.1 Level - 0 Data Flow Diagram (Context Diagram)

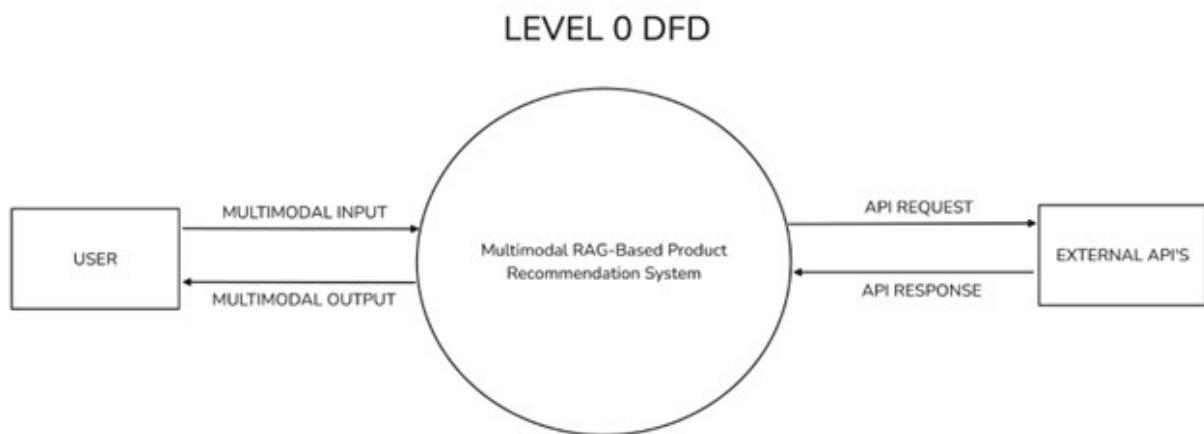


Figure 4: Level 0 DFD

14.2.2 Level - 1 Data Flow Diagram

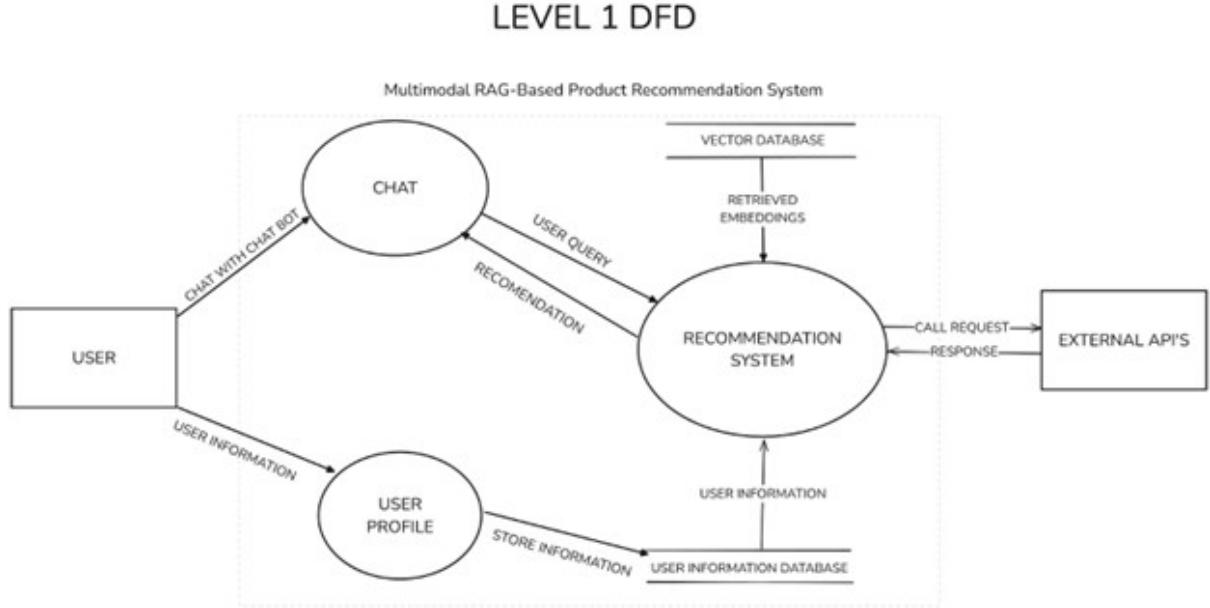


Figure 5: Level 1 DFD

15 CONCLUSION

The Multimodal Retrieval-Augmented Generation (RAG) system developed in this project represents a significant advancement in personalized recommendations for the fashion and cosmetics industries. By integrating textual descriptions, visual data, sustainability certifications, and real-time trend analysis, the system delivers accurate, user-centered, and environmentally conscious suggestions. Built on a scalable architecture using Flask, React, and vector databases, and powered by advanced AI services like OpenAI and Gemini, the platform ensures responsiveness and precision. This innovative system not only addresses critical gaps in traditional recommendation models but also promotes sustainable consumer behavior, setting a new benchmark for ethical and trend-aware AI applications in the industry.

System Design Document

16 Introduction

16.1 Purpose of this document

This Software Design Document (SDD) describes the architecture and system design of our project, Multimodal RAG-Based Product Recommendation System. The intended audience includes software developers, system architects, project managers, and stakeholders involved in our project. This document aims to provide a detailed and clear understanding of the platform's design.

The SDD covers the architectural design of our project Multimodal RAG-Based Product Recommendation System, providing a detailed look at the architecture. In addition to architectural details, this document includes activity diagrams and sequence diagrams of use cases to illustrate the dynamic aspects of the system. These diagrams help in understanding the flow of activities and interactions within the system, providing a clear visualization of how different components collaborate to achieve specific functionalities.

Furthermore, the SDD presents the UI design of our project, showcasing the visual and interactive elements that users will engage with. This ensures that stakeholders have a comprehensive view of the system's design, from backend architecture to user interface, facilitating better decision making and alignment throughout the development process.

16.2 Definitions, Acronyms, and Abbreviations

- **Multimodal RAG (Retrieval-Augmented Generation):** A system architecture that integrates retrieval mechanisms and generative models to process and combine multiple data modalities, such as text and images, to generate contextual outputs.
- **LLMs (Large Language Models):** Advanced machine learning models that are capable of understanding and generating human-like text based on large-scale datasets.
- **Vector Database:** A database for storing and retrieving high-dimensional data embeddings, often used in recommendation systems to compare and retrieve relevant results efficiently.
- **SDD:** Software Design Document, a document detailing the architectural and system design of the software.
- **UI/UX:** User Interface/User Experience, the overall experience of a user when interacting with a product, including its usability, accessibility, and aesthetics.
- **API (Application Programming Interface):** A set of protocols and tools that allow different software applications to communicate and exchange data.

- **Sustainability Filters:** Criteria applied to product data to identify and recommend eco-friendly items.
- **Embedding:** A mathematical representation of data (e.g., text, images) in a dense vector space, allowing the system to process multimodal data seamlessly.
- **Ethical AI Compliance:** The practice of ensuring the system adheres to guidelines minimizing bias and maintaining fairness, transparency, and user privacy.

17 System Overview

The Multimodal RAG-Based Product Recommendation System is designed to provide personalized recommendations in the fashion and cosmetics domains by leveraging state-of-the-art artificial intelligence technologies. This system integrates retrieval-augmented generation (RAG) architectures with multimodal data processing, combining user preferences, product metadata, and sustainability criteria to deliver tailored suggestions. It also supports dynamic trend analysis and real-time updates, ensuring the recommendations align with current market dynamics and user interests.

The system consists of several core components, including a user-friendly frontend built with React, a Flask-based backend for managing data processing and business logic, and a vector database for efficient embedding retrieval. External APIs, such as those provided by OpenAI, Gemini, and Hugging Face, are used to enrich the system with real-time trend data and advanced recommendation capabilities. The combination of text and image embeddings ensures a holistic understanding of user preferences and product attributes.

The architecture ensures high performance and scalability, allowing it to handle multiple user interactions and complex recommendation queries simultaneously. Users can explore recommendations through features like search and filter options, sustainability insights, and product comparisons, all of which are accessible via a seamless web interface. Furthermore, the system integrates ethical AI principles, such as bias mitigation and transparent decision-making, to enhance user trust and satisfaction.

By combining cutting-edge AI technologies, multimodal data handling, and user-centric design, this system aims to bridge the gap in sustainable and personalized recommendations, setting a benchmark in the fashion and cosmetics industry.

18 System Design

18.1 Architectural Design

The architecture of the Multimodal RAG-Based Product Recommendation System follows a layered three-tier design, ensuring modularity, scalability, and separation of concerns. The system is divided into three primary layers: Presentation Layer, Application

Layer, and Data Layer. Each layer has different roles and interacts with the other layers to provide a seamless user experience.

18.1.1 Layered Architecture

The layered architecture is a design pattern that separates the system into distinct, interconnected layers, where each layer is responsible for a specific aspect of the application. This architectural style promotes modularity, scalability, and maintainability by organizing components based on their functionality, allowing for independent development and testing.

In the context of our Multimodal RAG-Based Product Recommendation System, the layered architecture enables clear separation of user interactions, business logic, and data management. This separation ensures that changes made to one layer, do not disrupt the entire system, thereby reducing the risk of cascading failures.

We chose this architecture because it aligns with industry best practices for AI-driven web applications and effectively supports the processing and storage needs required for multimodal product recommendation systems. The clear division of responsibilities allows our development team to manage different system components simultaneously, increasing productivity and reducing development time.

18.1.2 Key Components

Our project's layered architecture consists of three primary layers, each fulfilling a distinct role within the system. These layers interact to deliver personalized recommendations, manage user data, and ensure the system operates efficiently. Components include:

- **Presentation Layer:** This layer manages user interaction and visualization. Built with React, it captures user inputs and displays real-time recommendations, product comparisons, and sustainability data.
- **Application Layer:** The core processing unit, implemented in Flask, handles business logic, API requests, and AI model integration. This layer processes user inputs, manages sessions, and communicates with external AI services to generate recommendations.
- **Data Layer:** Responsible for data storage and retrieval, this layer uses a database for structured data (user profiles, logs) and a vector database for vector embeddings. It ensures fast and scalable storage of product embeddings, allowing efficient search and recommendation generation.

18.1.3 Benefits of Layered Architecture

- **Modularity:** Each layer can be updated or replaced without affecting the entire system.

- **Scalability:** Layers can scale independently based on demand, ensuring optimal performance even during peak usage.
- **Security:** Sensitive user data is isolated in the Data Layer, enhancing security and preventing unauthorized access.
- **Maintainability:** Clear separation of concerns simplifies debugging and system enhancements.
- **Flexibility:** New features or external services can be integrated into specific layers, ensuring adaptability.

18.2 Decomposition Description

The Multimodal RAG-Based Product Recommendation System is designed using a modular three-tier architecture, with each layer handling distinct responsibilities. This decomposition ensures the system's scalability, maintainability, and flexibility. By dividing the system into functional layers Presentation Layer, Application Layer, and Data Layer the design allows for independent development and enhancement of each component, promoting efficient workflows and easier debugging.

18.2.1 Presentation Layer

The Presentation Layer acts as the user interface for the system, handling all interactions between the user and the application. It is responsible for capturing user inputs, displaying product recommendations, and visualizing real-time data, making it the most visible part of the architecture.

Responsibilities:

- Collects and processes user inputs such as login credentials, search queries, filter options, and profile settings.
- Displays interactive data, including product recommendations, trend insights, and sustainability details.
- Provides responsive and intuitive user interfaces that can adapt to different web browsers.
- Implements dynamic components using React, allowing for real-time updates without reloading the page.
- Facilitates secure communication with the backend via REST APIs for retrieving or submitting data.
- Ensures accessibility across browsers (Chrome, Safari, Firefox) for wider usability.

Key Features:

- **Authentication and Authorization:** Login and registration interfaces ensure secure access to personalized recommendations.
- **Recommendation Visualization:** Dynamic product recommendations are displayed with detailed trend analysis.
- **Search and Filter:** Users can search products using multi-criteria filters (e.g., price, sustainability).
- **Responsive Design:** Ensures seamless experiences across browsers and screen sizes.

18.2.2 Application Layer

The Application Layer is the core processing engine, managing the system's business logic and orchestrating data flow between the Presentation Layer and Data Layer. It handles AI model integration, user requests, and recommendation generation through advanced retrieval-augmented generation (RAG) techniques.

Responsibilities:

- Processes user requests, manages product searches and generates personalized recommendations.
- Hosts the Multimodal RAG framework, which combines text and image embeddings to enhance recommendation accuracy.
- Communicates with external AI services (e.g., OpenAI, Gemini, Hugging Face) to retrieve embeddings and analyze product trends.
- Manages user sessions, handles profile data, and updates product preferences based on real-time interactions.
- Implements Flask as the primary backend framework for creating and managing API endpoints.
- Facilitates interaction with external APIs to fetch data, analyze trends, and generate personalized recommendations.
- Incorporates error handling mechanisms to ensure system reliability and fallback solutions for API failures.

Key Features:

- **AI Model Integration:** Seamless integration with LLaMA, OpenAI CLIP, and Gemini for advanced product recommendations.
- **Real-Time Recommendations:** Continuously updates recommendations based on user interactions and external trend data.

- **Trend Monitoring:** Aggregates product trend data from APIs and processes it through multimodal embeddings.
- **Security Measures:** Implements user authentication and session management.

18.2.3 Data Layer

The Data Layer manages the storage, retrieval, and processing of all system data, including user profiles, embeddings, and product metadata. This layer plays a critical role in ensuring the system's responsiveness and scalability.

Responsibilities:

- Stores structured data such as user preferences, product information, and recommendation history.
- Manages high-dimensional embedding data for product recommendations using vector databases.
- Aggregates and stores trend data, sustainability certifications, and product details from external sources.
- Provides efficient data indexing and retrieval for fast recommendation generation.
- Maintains logs of user interactions to improve the accuracy of future recommendations.

Key Features:

- **Relational Database:** Stores essential user and product data, providing structured, fast retrieval.
- **Vector Database:** Manages embeddings for similarity searches and efficient product matching.
- **Metadata Storage:** Aggregates and stores external data, such as sustainability ratings, product reviews, and sales trends.
- **Data Caching:** Frequently accessed data is cached to reduce latency and improve response times.

18.2.4 External AI Services

To enhance the system's recommendation capabilities, the Application Layer integrates with external AI services that provide embedding generation, large language model (LLM) processing, and multimodal data analysis.

Services Possible to Use:

- **OpenAI (CLIP, GPT):** Generates embeddings for text and image data.
- **Gemini API (Google):** Provides real-time trend and sustainability data.
- **Hugging Face API:** Pre-trained models for NLP and sentiment analysis.
- **LLaMA (Meta AI):** Enhances recommendations through advanced LLM capabilities.

18.2.5 Communication Between Layers

The three layers communicate through REST APIs and database queries to facilitate data exchange. This modular design ensures that updates to one layer do not disrupt the others.

- **Presentation Layer to Application Layer:** RESTful API endpoints handle data requests and responses.
- **Application Layer to Data Layer:** SQL queries and embedding lookups manage interactions with databases.
- **Application Layer to External Services:** API calls to fetch embeddings and process product data.

18.3 System Modeling

18.3.1 Activity Diagrams

Use Case 2: Personalized and Interactive Recommendation Generation (Activity Diagram)

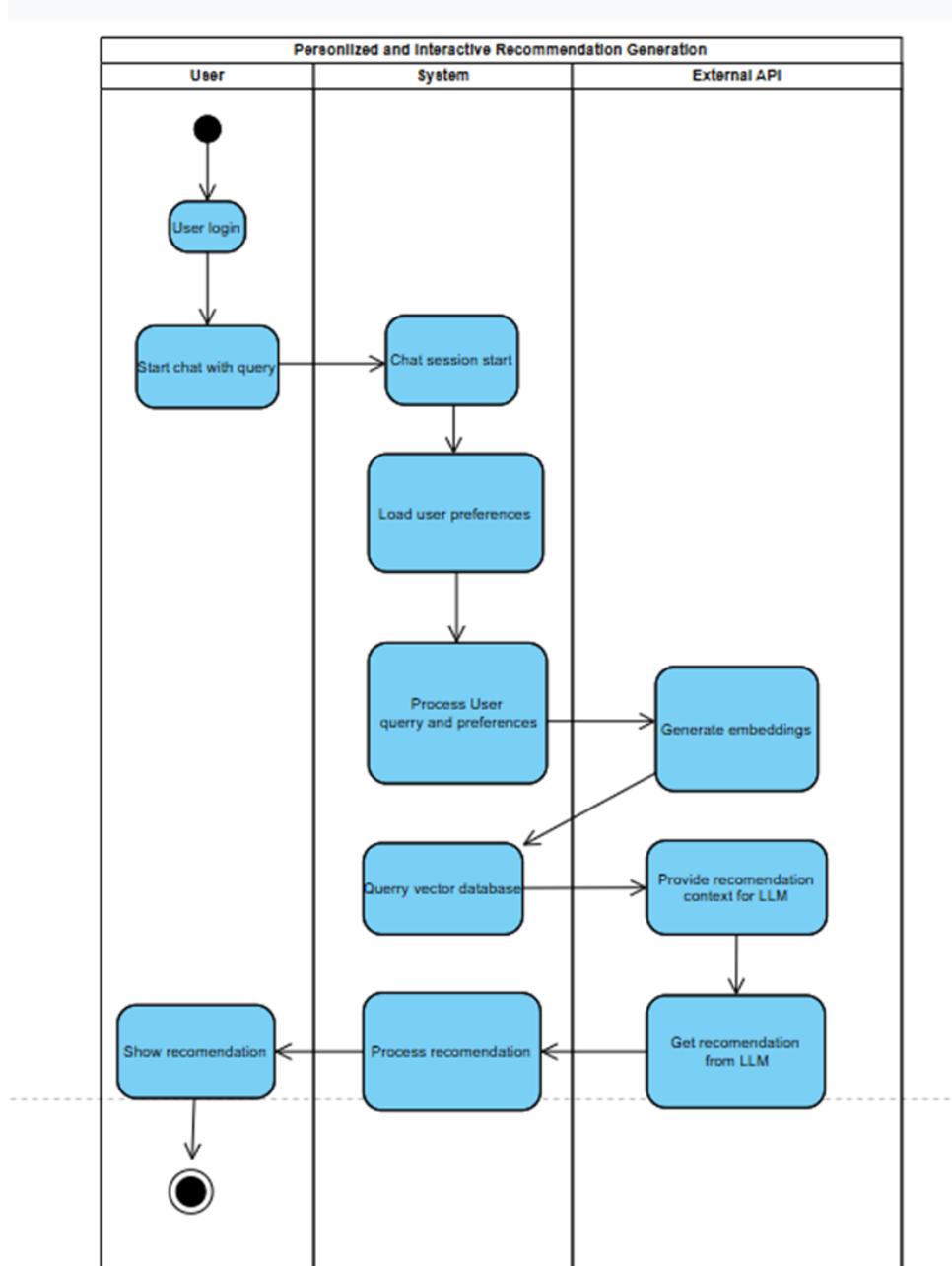


Figure 6: Use Case 2 Activity Diagram

Use Case 3: Trends Review (Activity Diagram)

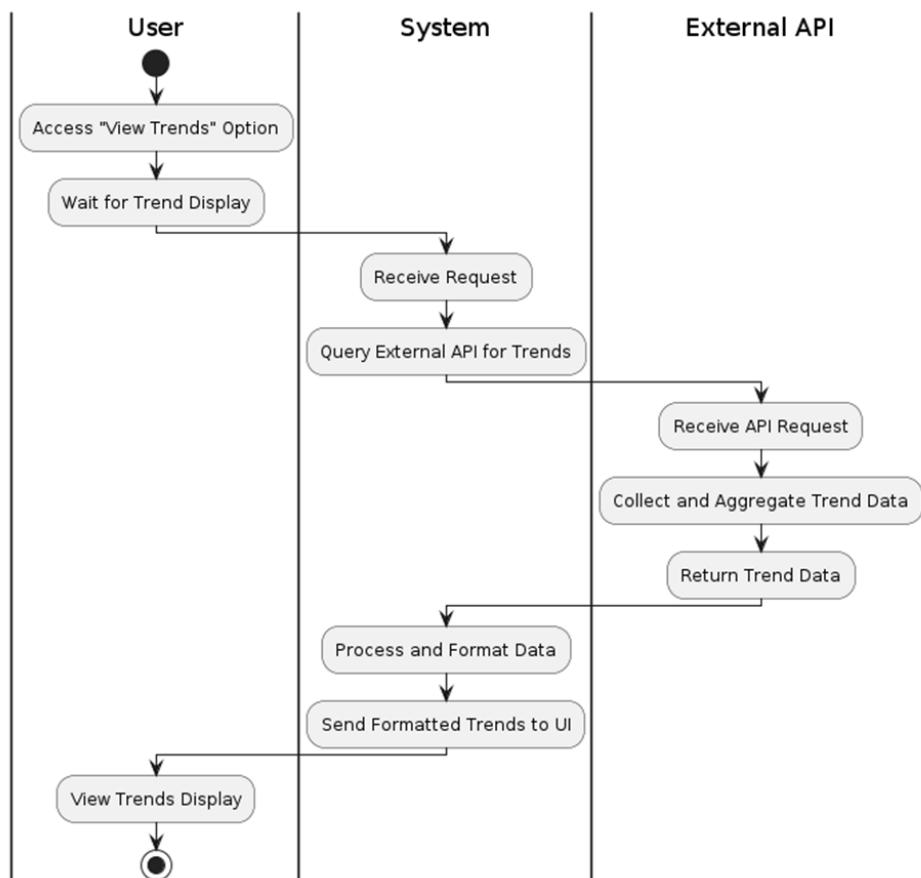


Figure 7: Use Case 3 Activity Diagram

Use Case 7: Search and Filter Options (Activity Diagram)

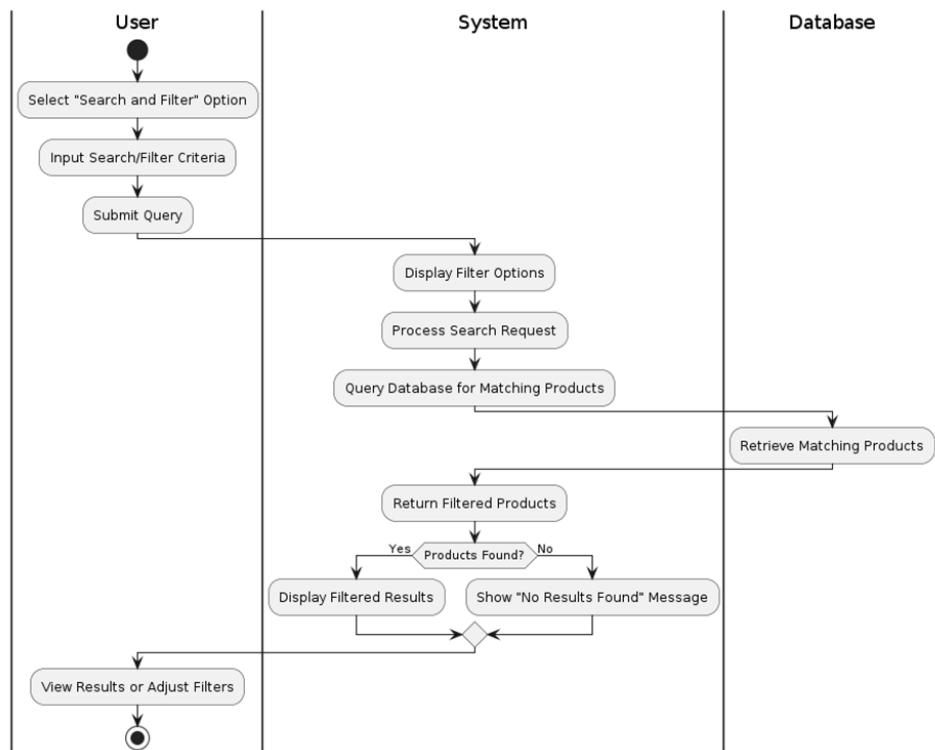


Figure 8: Use Case 7 Activity Diagram

Use Case 9: Product Comparison (Activity Diagram)

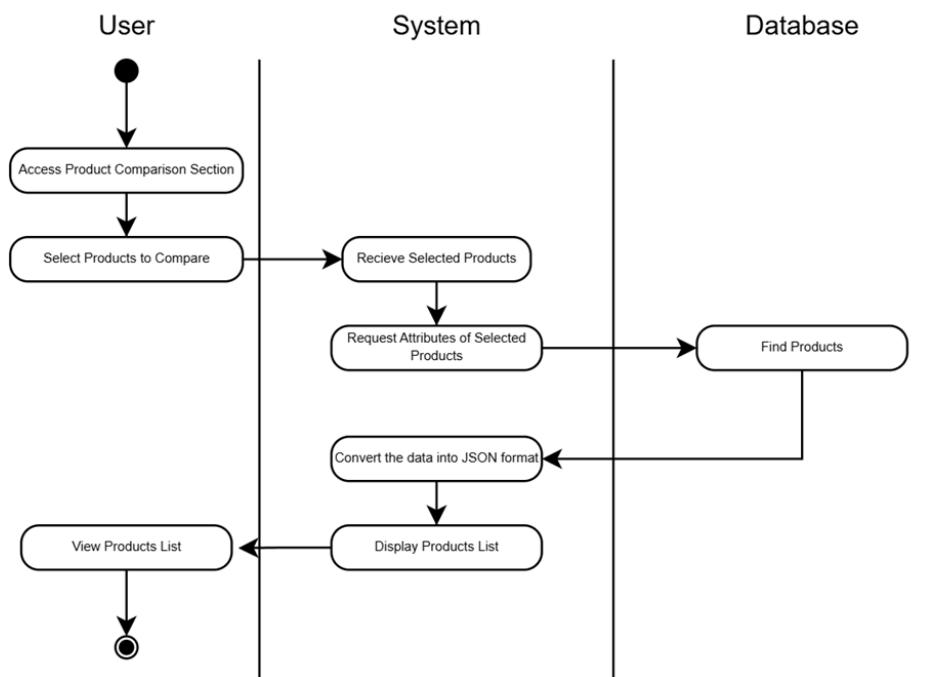


Figure 9: Use Case 9 Activity Diagram

18.3.2 Sequence Diagrams

Use Case 2: Personalized and Interactive Recommendation Generation (Sequence Diagram)

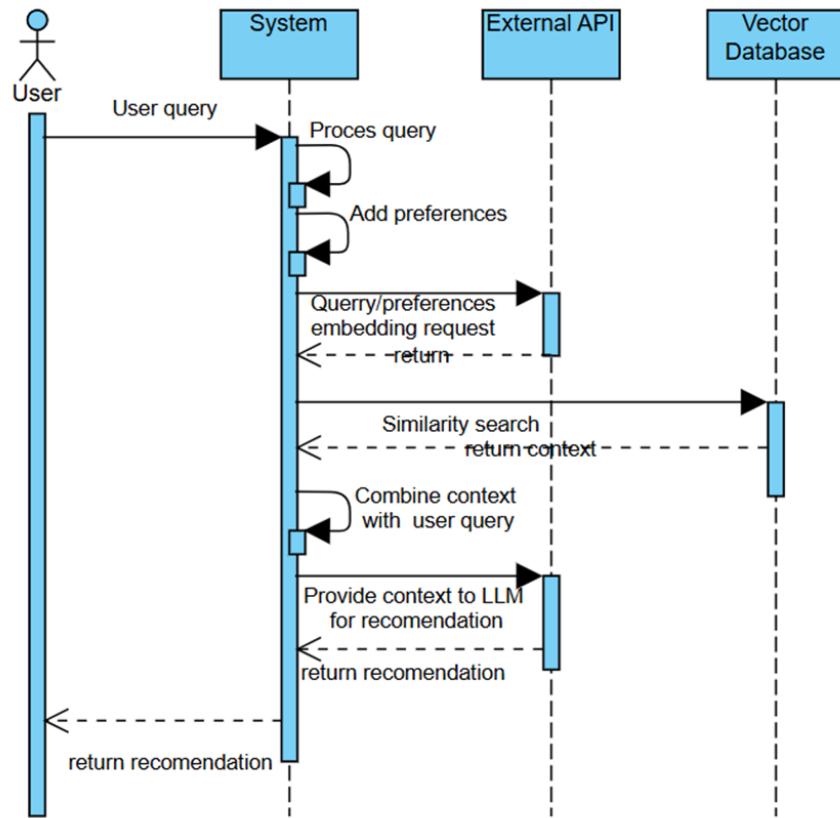


Figure 10: Use Case 2 Sequence Diagram

Use Case 3: Trends Review (Sequence Diagram)

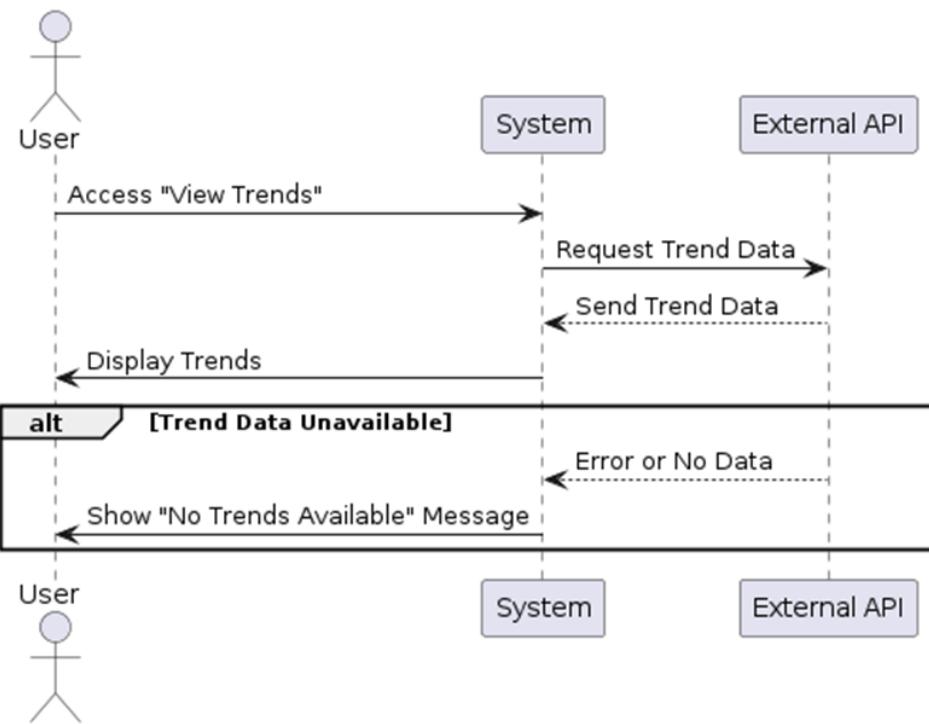


Figure 11: Use Case 3 Sequence Diagram

Use Case 7: Search and Filter Options (Sequence Diagram)

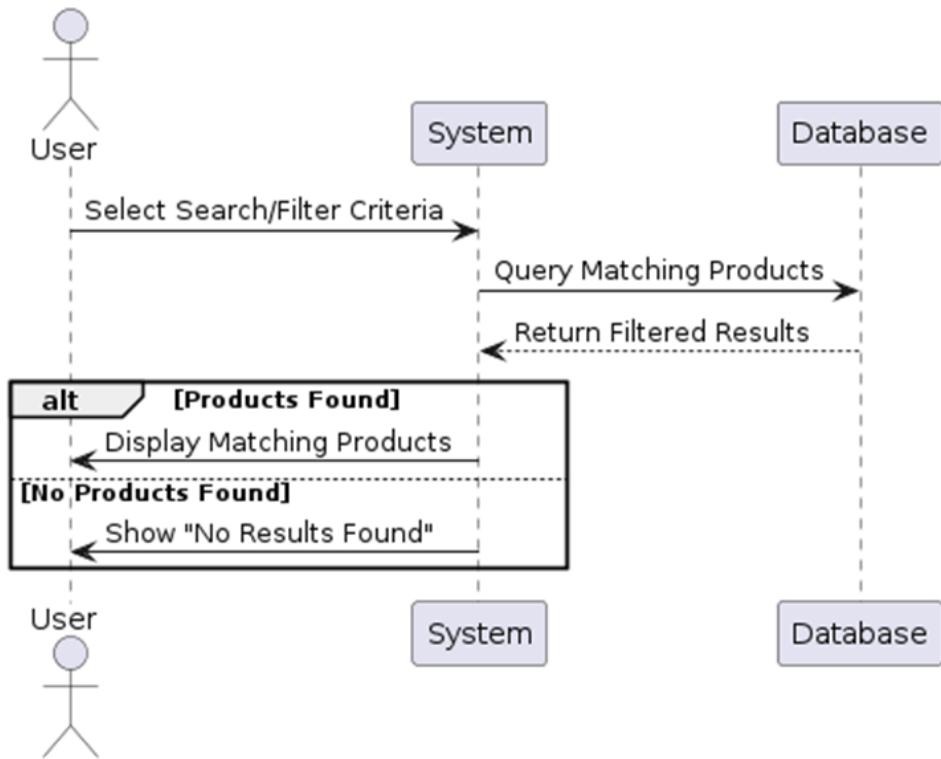


Figure 12: Use Case 7 Sequence Diagram

Use Case 9: Product Comparison (Sequence Diagram)

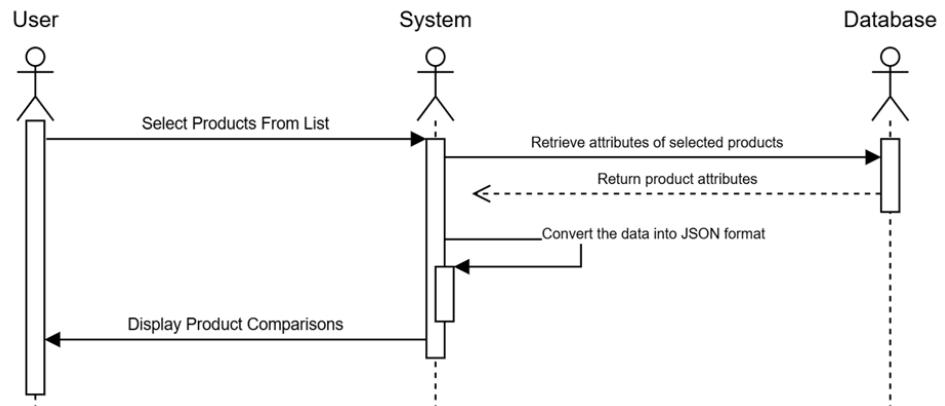


Figure 13: Use Case 9 Sequence Diagram

18.3.3 Class Diagram

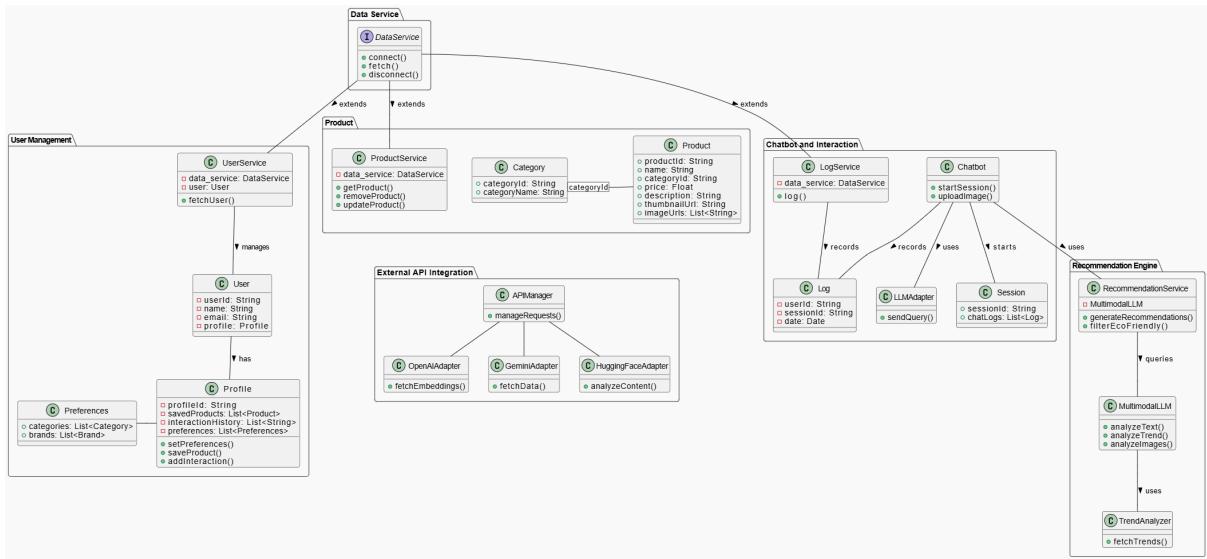


Figure 14: Class Diagram

19 User Interface Design

19.1 User Interface of Profile Page

The screenshot shows the 'My Profile' page of a shopping application named 'Glint'. At the top, there's a navigation bar with a logo, a search bar, and a user profile icon. Below the header, the main content area is titled 'My Profile'.

Favorites: This section displays three items: 'Neroli & Orchidee Eau de Toilette' (a clear bottle), 'L.12.12 Polo Shirt' (a close-up of a cream-colored polo shirt), and 'Rouge Volupte Shine Lipstick Balm' (a red lipstick).

Your favorite items: This section also lists the same three items with their prices: \$70.00 - \$140.00 for the perfume, \$90.00 for the shirt, and \$38.00 for the lipstick.

Your past searches: This section shows five circular thumbnails representing previous search queries: a perfume bottle, a makeup palette, a pair of high-heeled shoes, a man's polo shirt, and a handbag.

Your preferences: This section includes two buttons: 'Edit preferences' (with a note about preferred size, color, and brand) and 'How AI suggests items' (with a note about preferred size, color, and brand).

AI-Suggested Items: This section displays three new items: 'Orchidee Imprial Cream' (a white jar), 'L.12.12 Polo Shirt' (another view of the shirt), and 'Volupte Plump-in-Color Lip Balm' (three lip balms in different colors).

Footer: The footer contains links for 'About', 'Contact', 'Terms', and 'Privacy', along with social media icons for YouTube, Instagram, Facebook, and Twitter, and a copyright notice: '@2023 Glint.AI'.

Figure 15: Profile Page UI

19.2 User Interface of Chatbot Screen

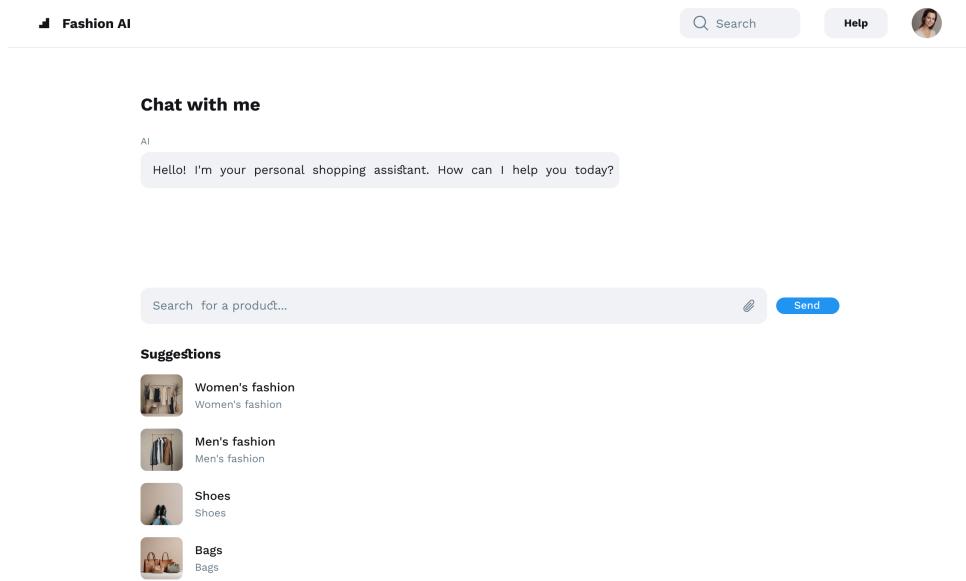


Figure 16: Chatbot UI

19.3 User Interface of Personal Recommendations Page

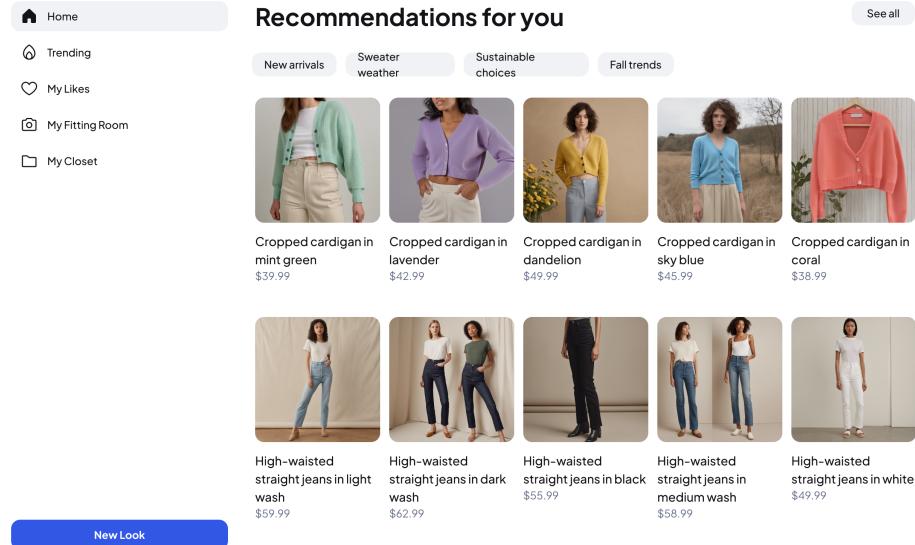


Figure 17: Personal Recommendations UI

19.4 User Interface of Home Page

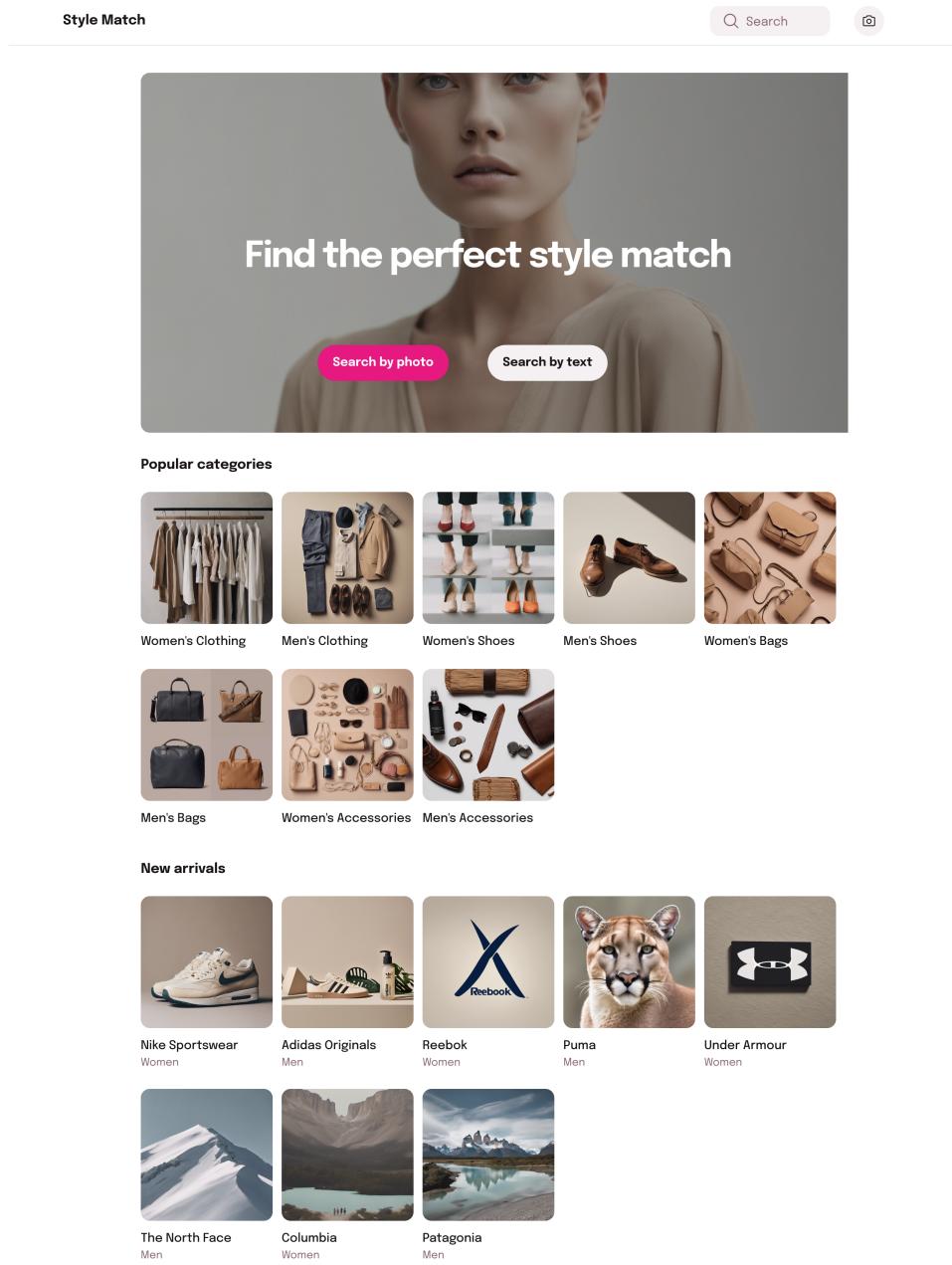


Figure 18: Home Page UI

19.5 User Interface of Login Page

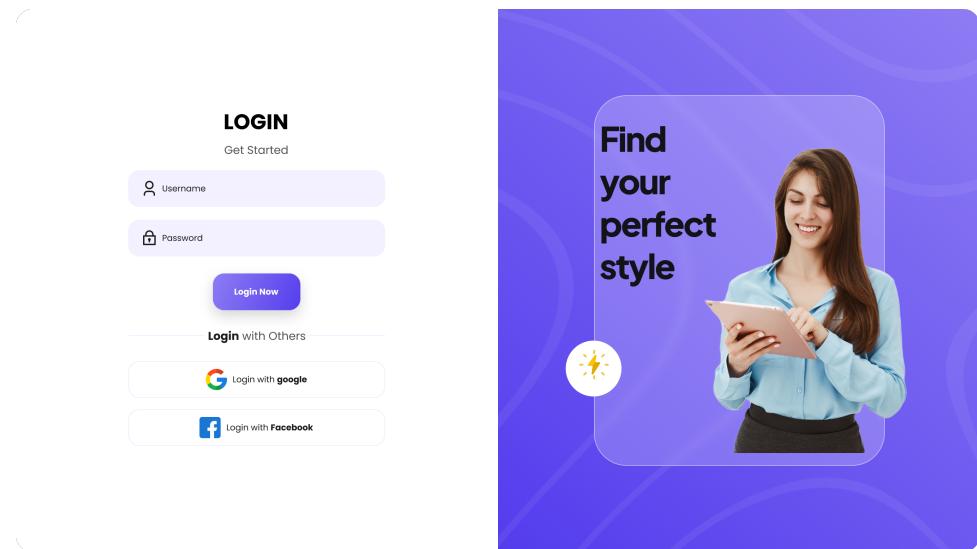


Figure 19: Login Page UI

19.6 User Interface of Trending Page

The screenshot shows the user interface of a fashion e-commerce website. At the top, there is a navigation bar with links for Women, Men, Home, Beauty, Sale, a Search button, and a Sign in button. Below the navigation bar is a search bar with a placeholder "Upload an image or type a query" and a "Search" button. To the left of the search bar is a large image of a woman wearing a purple sleeveless dress. To the right of the search bar is the text "Find your perfect style" and a smaller text "Upload an image or type a query". Below the search bar is a search input field with the text "ARA" and a "Search" button. Under the heading "Trending", there are ten items arranged in two rows of five. The items include a Dress (\$60), Sneakers (\$150), Earrings (\$30), Bag (\$200), Skincare (\$40), Pants (\$80), Sunglasses (\$100), Sandals (\$70), and a Jacket (\$120). Each item has a small image and its name and price below it. Below the trending section is a heading "Style Insights" followed by six circular images of women wearing different dresses.

Moda

Women Men Home Beauty Sale Search Sign in

Find your perfect style
Upload an image or type a query

ARA Search

Trending

Dress \$60	Sneakers \$150	Earrings \$30	Bag \$200	Skincare \$40
Pants \$80	Sunglasses \$100	Sandals \$70	Jacket \$120	

Style Insights

Figure 20: Trending Page UI

Conclusion

The development of this Multimodal Retrieval-Augmented Generation (RAG) system marks a significant step toward enhancing user experience in the fashion and cosmetics industries. By integrating important AI technologies and multimodal LLMs, the system not only personalizes recommendations but also promotes sustainable and trend-aware purchasing decisions. Users benefit from an interactive interface that seamlessly processes both text and image inputs, delivering highly relevant product insights.

The layered architecture ensures scalability, modularity, and easy maintenance, allowing for future improvements and integration of emerging AI models. Through features like sustainability filters and trend analysis, the system aligns with the growing demand for eco-conscious consumerism.

The accompanying SRS and SDD documents provide a detailed breakdown of the system's design, requirements, and functionality, serving as a foundation for continued development and refinement. This project sets the stage for broader adoption of AI technologies in the retail sector.

References

- [1] NVIDIA Developer. "An Easy Introduction to Multimodal Retrieval-Augmented Generation." *NVIDIA Developer Blog*, <https://developer.nvidia.com/blog/an-easy-introduction-to-multimodal-retrieval-augmented-generation/>. Accessed 11 Nov 2024.
- [2] Houlsby, N., Giurgiu, A., Jander, K., Schüpbach, M., Gerchow, M., & Senn, O. (2022). "Meta Learning with Memory-Augmented Neural Networks." *arXiv*, <https://arxiv.org/abs/2210.02928>.
- [3] Zhang, Y., et al. (2023). "Multimodal Retrieval-Augmented Generation (RAG) for AI Systems." *arXiv*, <https://arxiv.org/abs/2303.10868>.
- [4] Kim, J., Song, H., & Lee, S. (2022). "Generative Models with Efficient Transformers." *arXiv*, <https://arxiv.org/abs/2211.12561>.
- [5] Cheng, W., & Dong, L. (2024). "Deep Learning Techniques for Fashion Recommendation." *arXiv*, <https://arxiv.org/abs/2407.21439>.
- [6] Brown, J., & Lee, K. (2021). "Efficient NLP with Adaptive Models." *arXiv*, <https://arxiv.org/abs/2102.08871>.
- [7] Chen, L., & Yao, Z. (2024). "Advanced Techniques in Generative AI for Fashion." *arXiv*, <https://arxiv.org/abs/2408.08521>.
- [8] Zhang, Y., et al. (2023). "Multimodal Retrieval-Augmented Generation (RAG) for AI Systems." *arXiv*, <https://arxiv.org/pdf/2303.10868.pdf>.
- [9] Intelistyle. "Fashion Personalisation: The Ultimate Guide." *Inelistyle*, <https://www.inelistyle.com/fashion-personalisation-the-ultimate-guide/>. Accessed 08 Nov 2024.
- [10] PerfectCorp. "Top 10 AI-Powered Cosmetic Websites and Online Stores." *PerfectCorp Blog*, <https://www.perfectcorp.com/business/blog/makeup/top-10-ai-powered-cosmetic-websites-and-online-stores/>. Accessed 08 Nov 2024.
- [11] Chong, I. (2022). "AI-Powered Personalised Beauty Brands: Investing in AI and AR for Growth." *Cosmetics Design Asia*, <https://www.cosmeticsdesign-asia.com/Article/2022/09/08/ai-powered-personalised-beauty-brands-investing-in-ai-and-ar-are-growing-in-double-digits/>. Accessed 08 Nov 2024.
- [12] New York Post. "Formulate Personalized Shampoo Hair Care Review." *New York Post*, https://nypost.com/article/formulate-personalized-shampoo-hair-care-review/?utm_source=chatgpt.com/. Accessed 08 Nov 2024.
- [13] Intelistyle. "Top 12 Examples of Fashion Personalisation." *Inelistyle*, <https://www.inelistyle.com/top-12-examples-of-fashion-personalisation/>. Accessed 08 Nov 2024.
- [14] T. N. Prabhu, "Google Trends API," *Google Colab*, 2019. [Online]. Available: https://colab.research.google.com/github/Tanu-N-Prabhu/Python/blob/master/Google_Trends_API.ipynb. [Accessed: Dec. 4, 2024].

- [15] Stack Overflow, "What is the difference between functional and non-functional requirements?" *Stack Overflow*, 2013. [Online]. Available: <https://stackoverflow.com/questions/16475979/what-is-the-difference-between-functional-and-non-functional-requirements>. [Accessed: Dec. 4, 2024].
- [16] OpenAI, "CLIP: Connecting Text and Images," *OpenAI*, 2021. [Online]. Available: <https://openai.com/research/clip>. [Accessed: Dec. 4, 2024].
- [17] Pallets Projects, "Flask Documentation (3.1.x)," *Flask*, 2024. [Online]. Available: <https://flask.palletsprojects.com/en/stable/>. [Accessed: Dec. 4, 2024].
- [18] Google AI, "Gemini Models," *Google AI for Developers*, 2024. [Online]. Available: <https://ai.google.dev/gemini-api/docs/models/gemini>. [Accessed: Dec. 4, 2024].
- [19] Hugging Face, "The AI Community Building the Future," *Hugging Face*, 2024. [Online]. Available: <https://huggingface.co/>. [Accessed: Dec. 4, 2024].
- [20] Meta AI, "Llama 3.2," *Meta*, 2024. [Online]. Available: <https://www.llama.com/>. [Accessed: Dec. 4, 2024].
- [21] GeeksforGeeks. "Unified Modeling Language (UML) - Activity Diagrams." GeeksforGeeks, 2024. [Online]. Available: <https://www.geeksforgeeks.org/unified-modeling-language-uml-activity-diagrams/>. [Accessed: Dec. 20, 2024].
- [22] GeeksforGeeks. "Unified Modeling Language (UML) - Sequence Diagrams." GeeksforGeeks, 2024. [Online]. Available: <https://www.geeksforgeeks.org/unified-modeling-language-uml-sequence-diagrams/>. [Accessed: Dec. 20, 2024].
- [23] Figma. "UI Design Principles." Figma, 2024. [Online]. Available: <https://www.figma.com/resource-library/ui-design-principles/>. [Accessed: Dec. 23, 2024].
- [24] GeeksforGeeks. "Design Patterns - Architecture." GeeksforGeeks, 2024. [Online]. Available: <https://www.geeksforgeeks.org/design-patterns-architecture/>. [Accessed: Dec. 22, 2024].
- [25] P. Walpita. "Software Architecture Patterns: Layered Architecture." Medium, 2024. [Online]. Available: <https://priyalwalpita.medium.com/software-architecture-patterns-layered-architecture-a3b89b71a057>. [Accessed: Dec. 22, 2024].
- [26] Zeeshan. "The Weaknesses and Strengths of Layered Architecture in Software Development." Medium, 2024. [Online]. Available: <https://zeeshan01.medium.com/the-weaknesses-and-strengths-of-layered-architecture-in-software-development-81ba1206a17b>. [Accessed: Dec. 22, 2024].