# Technical Blueprint Report: Multimodal RAG for Personalized Food Recommendation

Divyansh Agarwal

## 1 Problem Context and Project Summary

### 1.1 The Challenge of Modern Recipe Discovery

The contemporary digital landscape offers an overwhelming abundance of culinary content, yet the process of recipe discovery remains fundamentally inefficient for many users. Existing food recommendation platforms and search engines typically operate on keyword matching and simplistic filtering mechanisms, such as cuisine type, dietary labels, or preparation time. While functional, this approach fails to capture the complex, multifaceted, and often subtle nature of a user's true intent. A user's desire for a meal is rarely a simple database query; it is a composite of dietary constraints (e.g., "low carbohydrate," "glutenfree"), taste preferences ("spicy but not overwhelmingly hot"), available ingredients ("what can I make with chicken, spinach, and rice?"), skill level ("a dish that is easy to make for a beginner"), and even aesthetic or social goals ("an impressive looking dessert for a dinner party").

This disparity between the richness of user intent and the shallow capabilities of conventional systems creates a significant satisfaction gap. Users are often forced to translate their nuanced desires into a series of rigid keywords, iteratively refining their search and manually sifting through irrelevant results. This process is time consuming, frustrating, and ultimately fails to deliver a truly personalized and helpful experience. The core problem is not a lack of content, but a lack of sophisticated understanding.

### 1.2 Proposed Solution: A Conversational, Multimodal Recommendation Engine

This document outlines the design for a conversational food recommendation engine engineered to bridge the aforementioned gap. The proposed solution moves beyond the paradigm of simple search to that of an intelligent culinary consultation. At its core, the system will leverage a state of the art Multimodal Retrieval Augmented Generation (RAG) architecture. This approach enables the system to understand and respond to user queries formulated in natural, conversational language.

The primary value proposition of this engine lies in its ability to jointly reason over multiple data modalities—specifically, the textual content of recipes (ingredients, instructions, nutritional information, user tags) and their corresponding images (presentation, color, texture, overall visual appeal). By processing and synthesizing information from both text and visuals, the system can deliver recommendations that are not only textually relevant but also contextually and aesthetically aligned with the user's request. For instance, it can interpret a query for a "vibrant, healthy summer salad" by retrieving recipes that are textually described as healthy and simultaneously feature images with bright, fresh looking ingredients.

The objective of this six week project is to develop a robust, functional proof of concept (PoC) that validates the efficacy of this multimodal RAG approach. The PoC will serve as the foundation for a new class of recommendation systems capable of a deeper, more human like understanding of user preferences in the culinary domain.

# 2 Dataset Specification and Preparation

## 2.1 Primary Dataset Selection: MealRec+

The successful implementation of a data driven system is contingent upon the quality and suitability of its underlying dataset. For this project, the `MealRec+` dataset, specifically the higher density `MealRecH+` variant, has been selected as the primary data source. This choice is strategic, as `MealRec+` provides a comprehensive, well structured collection of the exact modalities required for a sophisticated multimodal recommendation engine. While numerous other recommendation datasets exist, `MealRec+` offers a unique combination of features that make it exceptionally well suited for this project's objectives:

- **Rich Textual Data:** The dataset includes detailed textual information for each recipe, such as course names, ingredient lists, step by step cooking directions, user reviews, and descriptive tags. This textual richness is essential for building a deep semantic understanding of each dish.

- **Integrated Visual Data:** Each recipe entry is linked to an `image_url`, providing the critical visual modality needed for the multimodal architecture. Research indicates that users often rely on images first when selecting recipes, making this a crucial component for user preference modeling.

- **Valuable Metadata:** Crucially, `MealRec+` contains structured nutritional information and pre calculated healthiness scores based on Food Standards Agency (FSA) and World Health Organization (WHO) criteria. This data is a key strategic asset, enabling the system to address the growing user demand for health aware food recommendations and respond to complex, health conscious queries out of the box.

- **Relational Data:** The dataset also includes user course interaction data. While this collaborative filtering information will not be used in the initial RAG implementation, its presence provides a clear path for future enhancements, such as hybrid recommendation models or personalized fine tuning.

## 2.2 Data Schema, Modalities, and Preprocessing

A well defined preprocessing pipeline is required to transform the raw `MealRec+` data into a format suitable for embedding and ingestion. The primary input will be `course.csv`, which contains recipe level details. For each recipe, a unified document will be created by merging key textual fields: `coursename`, `ingredients`, `cookingdirections`, `tags`, and `nutritions` into one coherent text block to ensure comprehensive semantic coverage. The `imageurl` field will be used to fetch the corresponding image for multimodal processing.

A Python based pipeline will automate this process using libraries such as Pandas for data manipulation and Requests for image retrieval. The steps include:

1. Parsing the `MealRec+` dataset files, focusing on `course.csv` and the health score files.

2. Constructing unified text documents and appending nutritional and health data.

3. Downloading associated images with error handling for missing or invalid URLs.

4. Storing the processed text and image files in an Amazon S3 bucket as structured JSON objects, forming the canonical data source for the embedding stage.
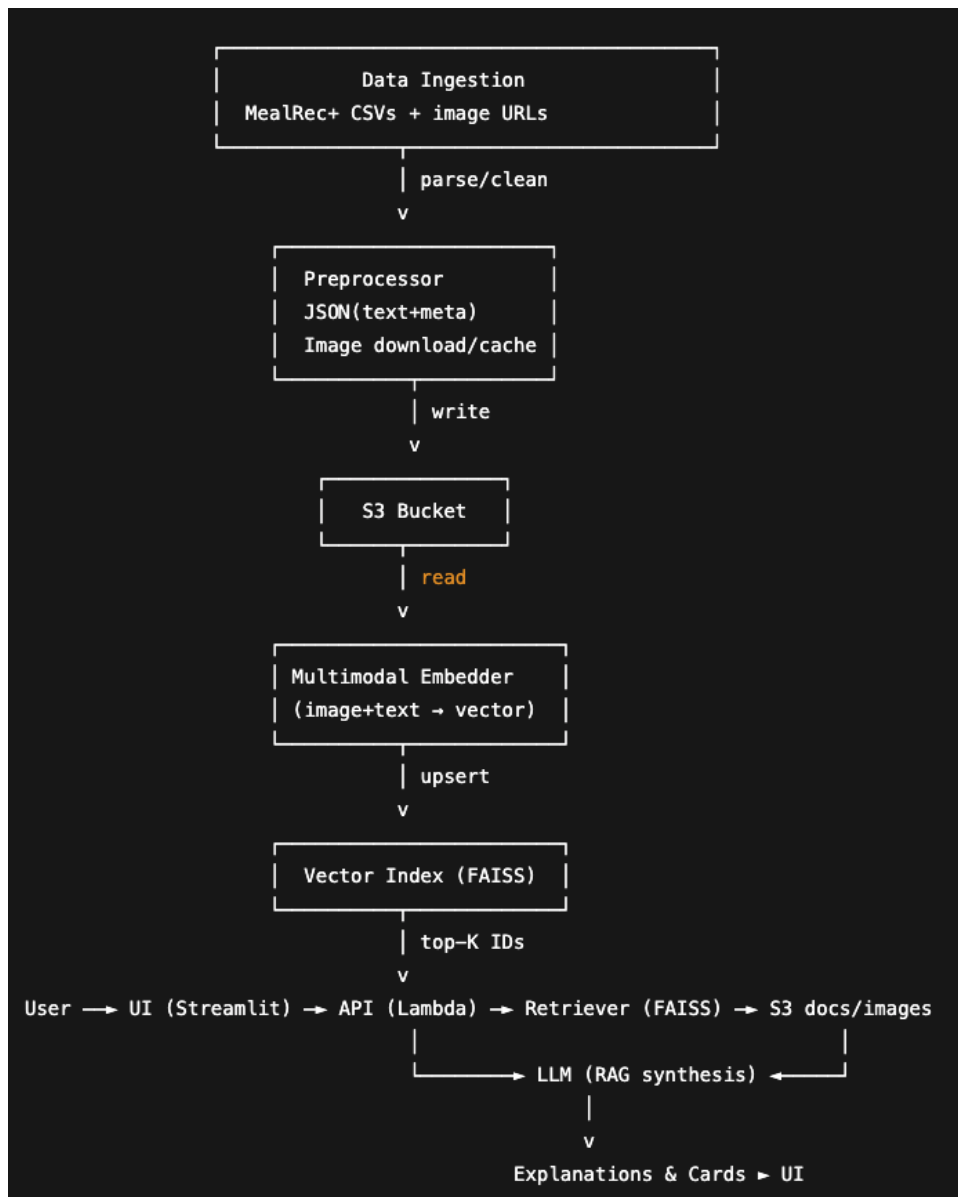
# 3 Planned System Architecture

The proposed system integrates two main workflows:

- **Data Ingestion and Embedding (Writer Path):** Responsible for preparing, embedding, and storing recipes in a multimodal vector index.

- **Query and Retrieval (Reader Path):** Handles real time user queries, retrieves the most relevant recipes, and generates personalized responses through a conversational interface.

This architecture blends managed cloud services (for storage, scalability, and model inference) with optimized compute for fast retrieval, ensuring both performance and maintainability.

## 3.1 Data Flow Blueprint



```
        ┌──────────────────────────────────────┐
        │            Data Ingestion            │
        │       MealRec+ CSVs + image URLs     │
        └──────────────────────────────────────┘
                         │ parse/clean
                         v
              ┌────────────────────────┐
              │  Preprocessor          │
              │  JSON(text+meta)        │
              │  Image download/cache   │
              └────────────────────────┘
                         │ write
                         v
                 ┌──────────────────┐
                 │   S3 Bucket      │
                 └──────────────────┘
                         │ read
                         v
              ┌────────────────────────┐
              │ Multimodal Embedder    │
              │ (image+text → vector)   │
              └────────────────────────┘
                         │ upsert
                         v
              ┌────────────────────────┐
              │ Vector Index (FAISS)   │
              └────────────────────────┘
                         │ top-K IDs
                         v
User ──► UI (Streamlit) ──► API (Lambda) ──► Retriever (FAISS) ──► S3 docs/images
                         │                                          │
                         └──────────────► LLM (RAG synthesis) ◄─────┘
                                                │
                                                v
                                  Explanations & Cards ► UI
```

## 3.2 Subsystem Deep Dive

### 3.2.1 Data Ingestion and Embedding Pipeline (The "Writer" Path)

This pipeline processes the source data and populates the vector store. It runs as a one time batch job for the initial dataset but can be adapted to an event driven model for future updates.

- **Source Data:** Begins with preprocessed recipe text (JSON) and images stored in an **Amazon S3** bucket for durability and easy integration.

- **Embedding Model:** Uses a multimodal embedding model on **Amazon Bedrock** to jointly encode image and text features into a shared semantic space. This "direct multi-modal embedding" approach preserves visual details without relying on text only captions.

- **Orchestration:** A Python script on an **EC2 instance** iterates through S3 data, generates embeddings, and writes them to the vector index. For incremental updates, an **AWS Lambda** can be triggered by S3 upload events.

- **Output:** The final output is a collection of dense vectors, each representing the full semantic meaning of a recipe, stored for retrieval.

### 3.2.2 Vector Storage and Retrieval Layer

This layer stores high dimensional embeddings and performs efficient similarity searches for user queries.

- **Vector Store:** Implements **FAISS (Facebook AI Similarity Search)** for high performance approximate nearest neighbor retrieval, hosted on a memory optimized **EC2 instance**. The index is loaded fully into RAM for sub second response times.

- **Index Type:** Uses `IndexIVFPQ`, which partitions and compresses vectors for efficient memory use while maintaining accuracy.

- **Metadata Mapping:** A lightweight mapping between FAISS IDs and recipe IDs is maintained using an in memory dictionary or **Redis** store.

### 3.2.3 Generative Core and Orchestration (The "Reader" Path)

This workflow handles user queries, retrieves relevant content, and generates conversational responses.

- **API Layer:** A managed **API Gateway** exposes a secure endpoint that triggers a central **AWS Lambda** function.

- **Workflow:** The Lambda function embeds the user query, connects to the FAISS server, retrieves top K matching recipes, and fetches metadata from S3. It then constructs a structured prompt with the user query and retrieved recipes.

- **Generative Model:** The prompt is sent to a large language model via **Amazon Bedrock** (e.g., Claude 3 Sonnet) to synthesize a personalized response grounded in the retrieved recipes.

- **Response Handling:** The generated text and recipe data are formatted as JSON and returned to the frontend for display.

# 4 User Interface (UI) Plan

## 4.1 Technology Choice: Streamlit

The user interface for this proof of concept will be developed using **Streamlit**, a Python based framework that enables the rapid creation of interactive, data driven web applications. This choice allows for fast prototyping and easy integration with the backend without requiring extensive frontend development. Building directly in Python keeps the focus on the AI functionality and accelerates iteration, which is ideal for a six week development timeline.

## 4.2 User Interaction Flow and Key Components

The UI will prioritize a simple, conversational experience that makes interactions intuitive and visually engaging. It will include two main interactive components:

- **Conversational Chat Interface:**

  The chat interface will serve as the primary mode of interaction. A text input box at the bottom of the screen will let users type natural language queries (e.g., "Show me a healthy vegetarian dinner under 600 calories"). Above this, a scrollable chat panel will display the dialogue history, preserving conversational context.

- **Recommendation Display Area:**

  After each query, the AI's response will appear in the chat history, followed by a set of visually rich *recipe cards*. Each card will contain:

  - **Recipe Image** – the key visual cue for the dish.
  - **Title** – clickable recipe name.
  - **AI Generated Rationale** – a short explanation (1–2 sentences) on why this dish fits the query.
  - **Tags** – keywords such as "Vegan," "Quick Meal," or "High Protein."

This layout enhances user trust and transparency by showing not only recommendations but also the reasoning behind them. Each recommendation includes a brief explanation generated by the model ("why this recipe was selected"), allowing users to understand the logic of the system. The result is a personalized, conversational, and visually grounded experience that clearly demonstrates the benefits of multimodal retrieval.

# 5 Innovation and Anticipated Challenges

## 5.1 Core Innovation: Cross Modal Semantic Understanding

The core innovation of this project lies in implementing a direct multimodal embedding RAG pipeline for a real world recommendation task. Traditional systems are either unimodal or rely on simplified multimodal approaches that convert images into textual captions, losing much of the visual richness and context.

In contrast, this system jointly embeds both text and images into a shared semantic space, enabling cross modal retrieval. It can match user intent expressed in text (e.g., "light and vibrant dish perfect for summer") with visually relevant recipes—images characterized by bright colors, fresh ingredients, and aesthetic composition—even if those exact words do not appear in the text. This ability to reason across modalities allows the system to interpret abstract or stylistic preferences and make more human like, visually aligned recommendations.

## 5.2  Technical Challenges and Mitigation Plan

Developing a multimodal RAG pipeline involves several integration and scalability challenges. The table below summarizes key challenges and their planned mitigation strategies.

| Challenge | Description | Mitigation Strategy |
|---|---|---|
| **Semantic Misalignment** | Text and image embeddings may not align for certain culinary concepts, causing irrelevant results. | Use a large pre trained multimodal model for better alignment; validate with "challenge queries" during testing. |
| **Retrieval Latency** | Query embedding and FAISS search could exceed acceptable response time. | Host FAISS on a memory optimized EC2 instance, preload index in RAM, and use an efficient index type (`IVFPQ`). |
| **Model Hallucination** | The LLM may generate inaccurate or fabricated recipe details. | Apply strict prompt constraints and display retrieved sources for user transparency. |
| **Compute Cost** | Embedding and hosting FAISS continuously may increase cloud expenses. | Start with smaller EC2 instances; monitor usage and scale only if needed. |

# 6  Six Week Implementation Timeline

## 6.1  Phased Development Plan

The project will follow an agile, six week plan with progressive integration of key components—starting from data setup to a fully functional, tested prototype.

| Week | Focus | Key Tasks | Deliverables |
|---|---|---|---|
| 1 | Environment Setup & Data Ingestion | Configure AWS (IAM, S3, VPC). Parse and clean the MealRec+ dataset, store JSON & image files in S3. | Structured dataset hosted in S3. |
| 2 | Embedding Pipeline & Vector Index | Implement batch embedding via Bedrock Titan, build FAISS index on EC2, test data retrieval. | Working FAISS index ready for search. |
| 3 | Backend RAG API | Develop Lambda orchestration for query embedding, FAISS retrieval, and Claude model call. Deploy via API Gateway. | Working API returning contextual responses. |
| 4 | Streamlit UI Integration | Build chat based interface. Connect to API, render query responses. | Interactive prototype with text based output. |
| 5 | Testing & Optimization | Add recipe cards, improve prompts, run benchmark queries, measure latency, refine results. | End to end tested RAG system with improved UX. |
| 6 | Final Deployment & Docs | Debug, clean codebase, add README/setup docs, and rehearse demo. | Stable, documented system ready for presentation. |

# 7 Responsible AI Reflection

Developing intelligent systems demands careful consideration of ethics, fairness, and user safety. This project follows Responsible AI principles, with proactive steps to ensure fairness, transparency, safety, and privacy throughout the system lifecycle.

## 7.1 Algorithmic Bias and Fairness

- **Issue:** Real world datasets often reflect cultural or demographic biases. The `MealRec+` dataset may overrepresent Western cuisines, leading the recommender to favor these recipes.

- **Mitigation:** Although full debiasing is outside the current project scope, evaluation will include fairness testing by querying recipes from diverse cuisines. For future iterations, data augmentation or diversity based re ranking can balance recommendations.

## 7.2 Transparency and Explainability

- **Approach:** The RAG architecture ensures transparency by grounding responses in retrieved recipe documents.

- **Implementation:** The Streamlit UI explicitly shows which recipes informed each recommendation, along with an AI generated rationale explaining *why* a dish was suggested. This supports user trust and interpretability.

## 7.3 Safety and Reliability

- **Concern:** Food recommendations may pose health risks if they contain allergens or promote unhealthy eating habits.

- **Mitigation:** Safety will be reinforced through:

  1. **Prompt Guardrails:** The LLM will be instructed to avoid medical claims and include a disclaimer such as, "Consult a nutritionist for specific dietary advice."
  2. **Curated Data Source:** The system uses only the vetted `MealRec+` dataset to avoid misinformation or unsafe content.

## 7.4 Data Privacy and Security

- **Principle:** The system follows a privacy by design approach.

- **Implementation:** User queries are processed in memory via AWS Lambda and not stored. All communications use HTTPS/TLS encryption. Access to AWS services (S3, EC2, Bedrock) is tightly restricted under least privilege IAM roles. A future production deployment would include anonymization and data retention policies compliant with privacy standards.