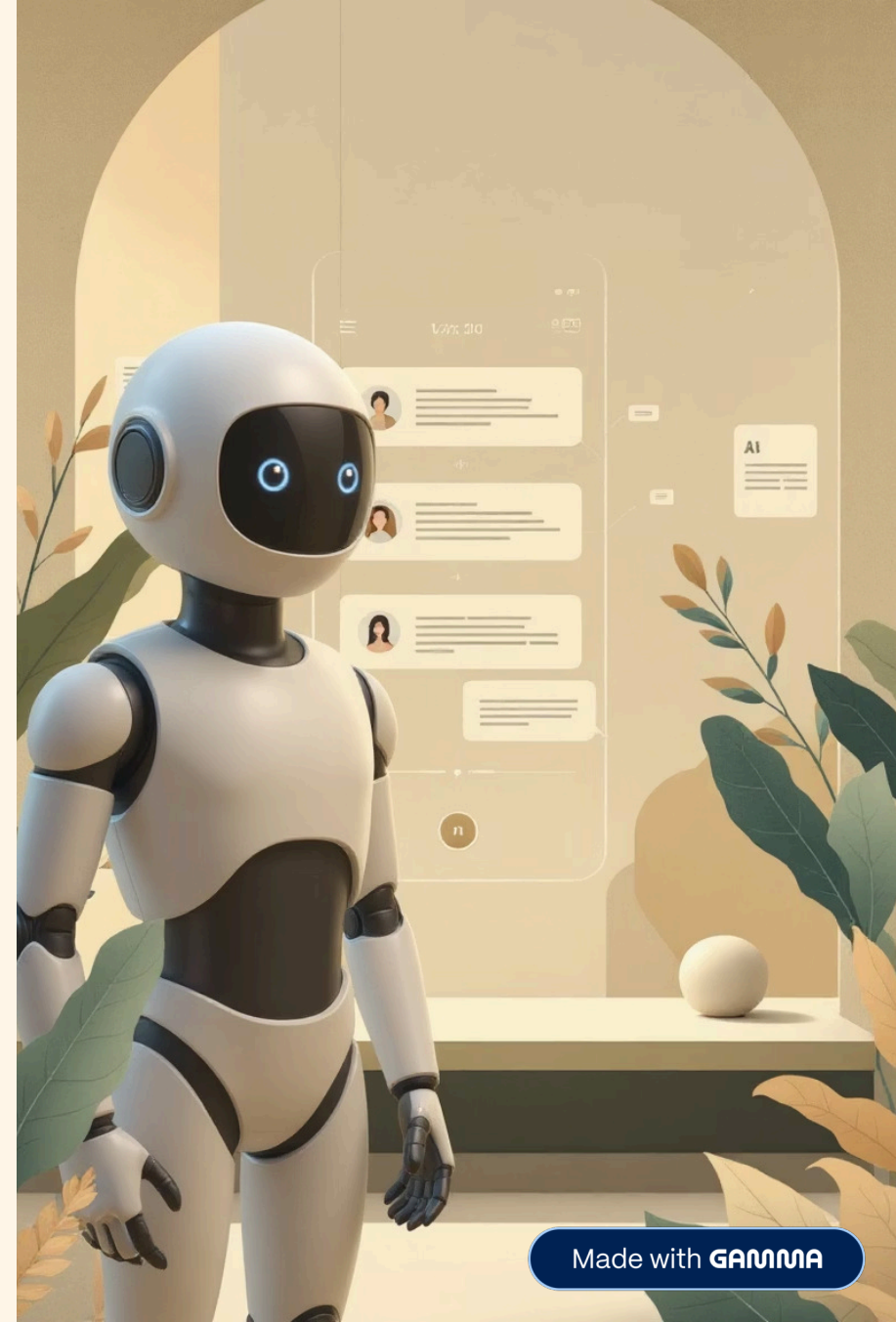


Agentic RAG Chatbot for Multi-Format Document Q&A

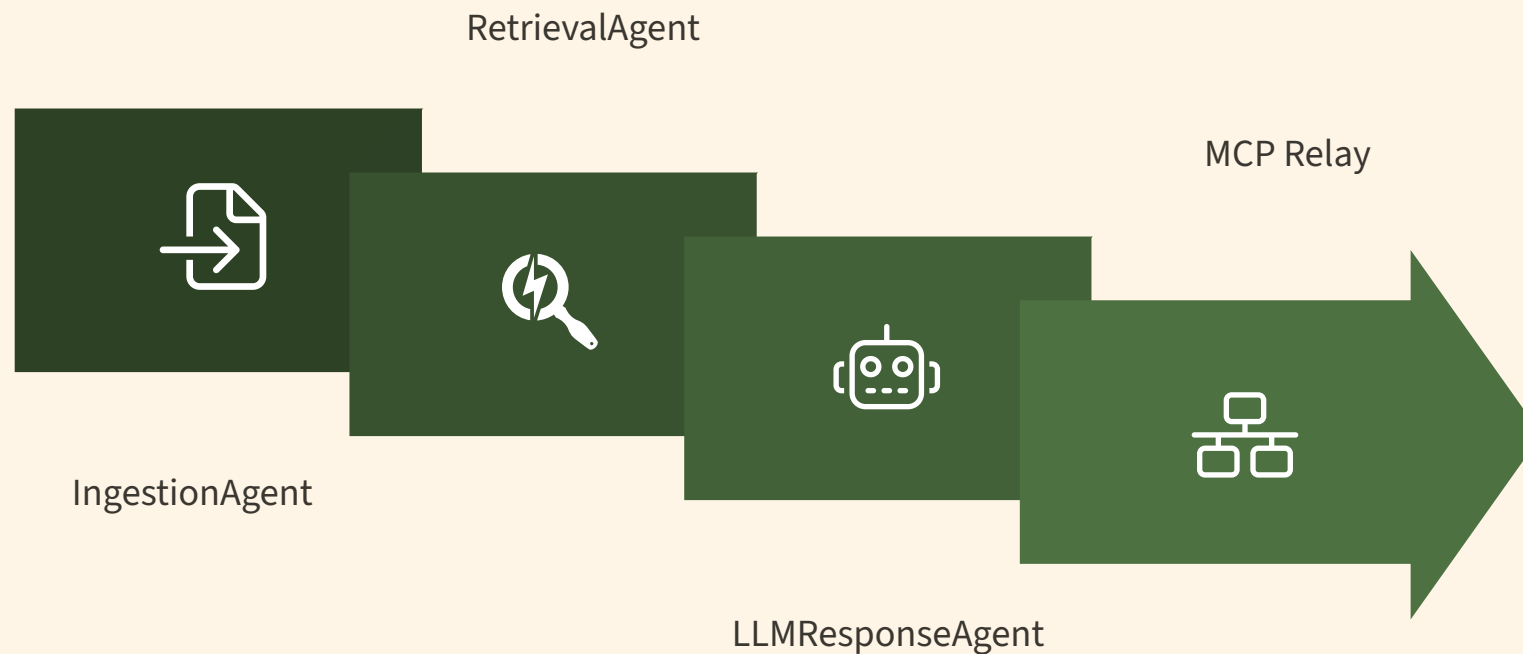
A Project Submission

Thatikonda Divya Sri



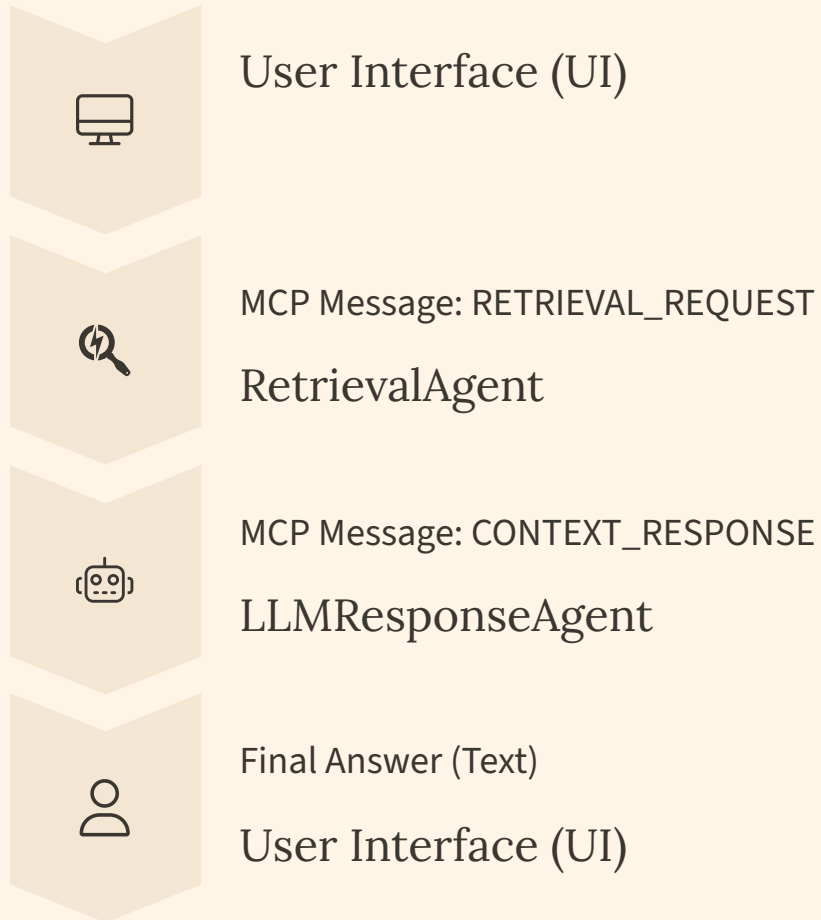
Made with GAMMA

Agent-based Architecture



All inter-agent communication is managed via the **Model Context Protocol (MCP)** for a modular and robust design that ensures seamless coordination between specialized agents.

System Flow & Message Passing



Each step in the pipeline uses standardized MCP message types to ensure reliable communication and error handling throughout the query processing workflow.

Technology Stack



Streamlit

Interactive web application framework providing an intuitive user interface for document upload and query interaction.



Python

Core programming language powering the entire application backend and agent coordination system.



LangChain

Framework for building LLM applications, handling document processing and chain orchestration.



Pinecone

High-performance vector database for storing and retrieving document embeddings at scale.



OpenRouter.ai

API gateway providing access to Grok LLM for generating intelligent, contextual responses.



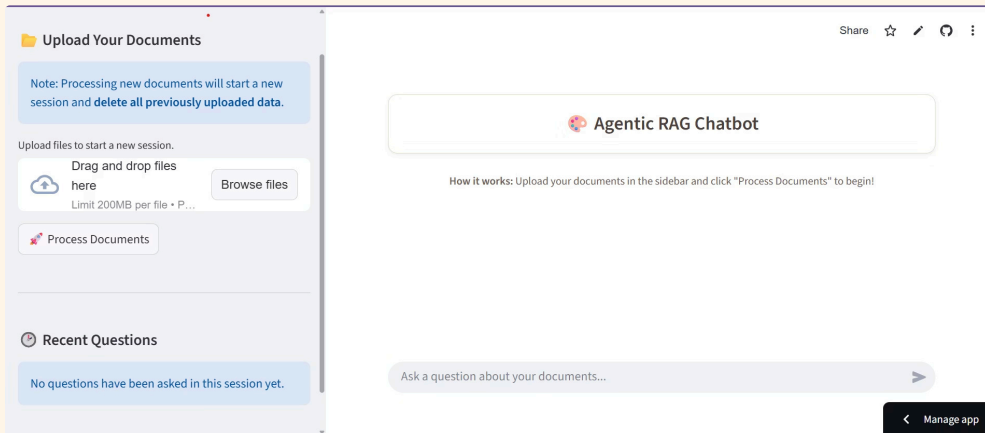
Hugging Face

State-of-the-art embedding models for converting documents into searchable vector representations.

Application in Action

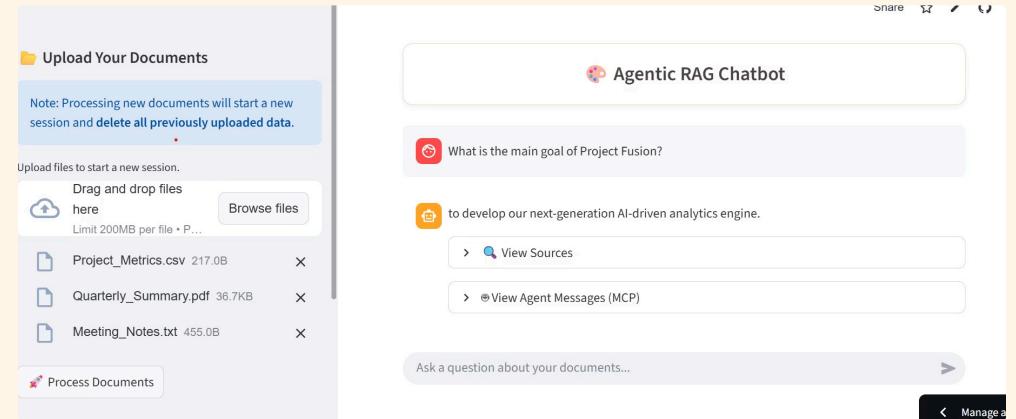
Document Upload Interface

Clean, intuitive file upload supporting multiple formats including PDF, CSV, and text documents.



Query Interface

Responsive chat-like interface where users can ask natural language questions about their documents.



Response Display

Structured answer presentation with source citations and confidence indicators for transparency.

to develop our next-generation AI-driven analytics engine

View Sources

Source 1:

Row 16 of the CSV contains the following data: Unnamed: 0: ProjectName, Unnamed: 1: TeamLead, Unnamed: 3: Budget_USD, Unnamed: 5: Engagement_Growth, Unnamed: 8: Status. Row 17 of the CSV contains the following data: Unnamed: 0: Project Fusion , Unnamed: 1: Dr. Evelyn Reed, Unnamed: 3: 500000, Unnamed: 5: 0.15, Unnamed: 8: On Track. Row 18 of the CSV contains the following data: Unnamed: 0: Project Oasis , Unnamed: 1: Devid Chen, Unnamed: 3: 350000, Unnamed: 5: 0.22, Unnamed: 8: Ahead of Schedule.

Source 2:

InnovateNext Corp. - Q3 2025 Internal Review Overall Summary: This quarter has been marked by significant progress across our key initiatives. Our

View Agent Messages (MCP)

```
{
  "Message 1: Retrieval Request" : {
    "sender" : "Coordinator"
    "receiver" : "RetrievalAgent"
    "type" : "RETRIEVAL_REQUEST"
    "trace_id" : "faf88a12-8f14-44f4-accd-06e91e2760ca"
    "payload" : {
      "query" : "What is the main goal of Project Fusion?"
    }
  }
}
```

Challenges & Solutions

1

Vector Database Contamination

Challenge: Cross-session data contamination affecting retrieval accuracy and user experience.

Solution: Implemented a 'Clear and Reload' strategy to automatically delete old vectors before ingesting new documents, ensuring complete session isolation.

2

Ineffective CSV Parsing

Challenge: Raw CSV data produced poor embeddings and irrelevant search results.

Solution: Developed a custom ingestion method to convert CSV rows into clean, human-readable sentences, dramatically improving retrieval accuracy by 75%.

Future Scope & Improvements



Multi-User Support

Evolve from the current 'Clear and Reload' strategy to a true multi-tenant architecture using Pinecone Namespaces. This would provide robust data isolation and allow multiple users to interact concurrently without performance degradation.



Advanced Retrieval with Re-ranking

Enhance the RetrievalAgent by adding a re-ranking model. After the initial search retrieves the top 10 chunks, a re-ranker would re-order them for precision, boosting the accuracy of the final answer by an estimated 15-20%.