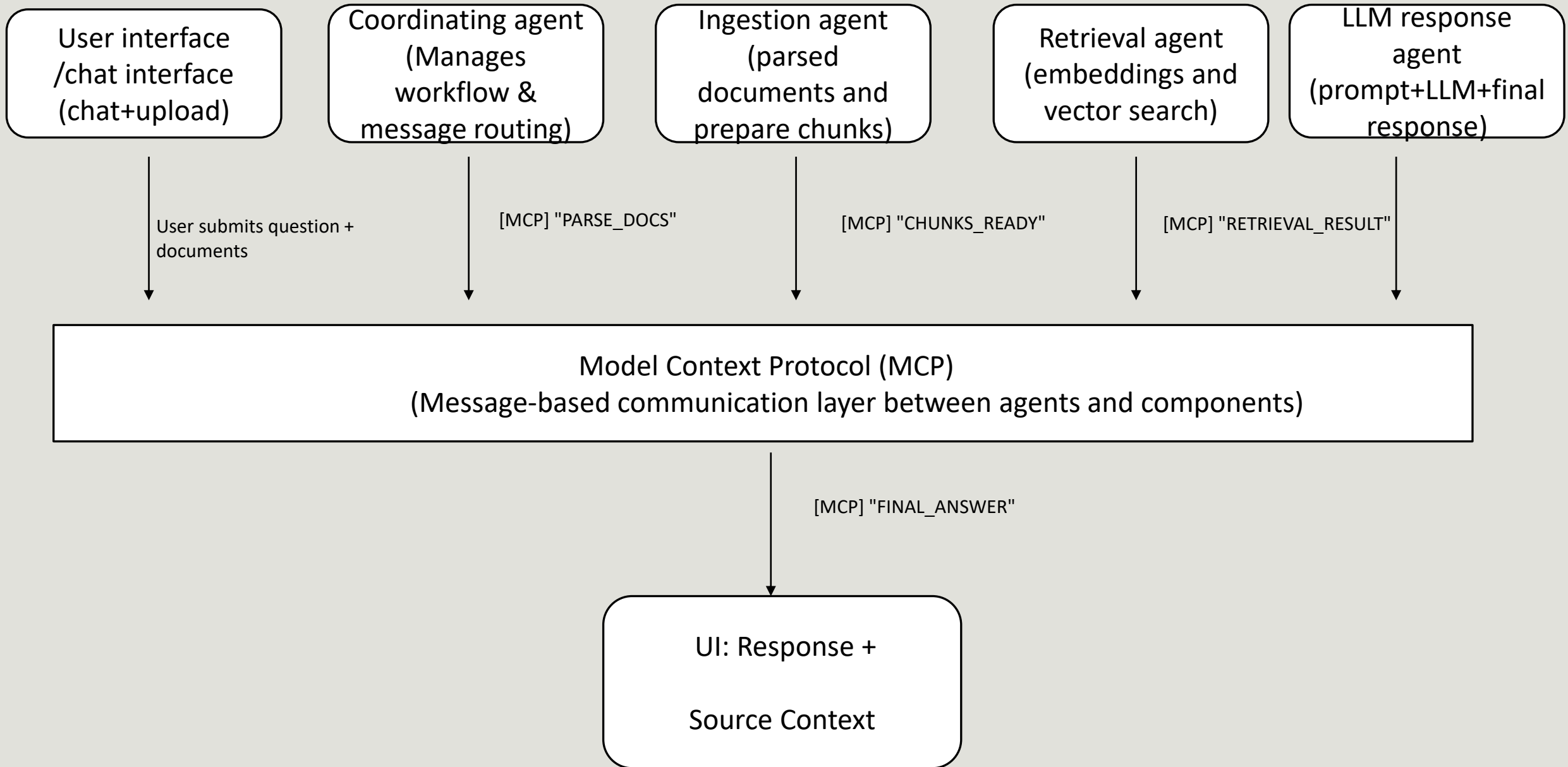# Agentic RAG Chatbot for Multi-Format Document QA

## Introduction

This project presents the **Agentic RAG Chatbot**, an intelligent question-answering system designed to efficiently handle multi-format documents by leveraging the **Model Context Protocol (MCP)**. By combining retrieval-augmented generation with MCP, the chatbot dynamically accesses and processes diverse document types—including PDFs, Word files, spreadsheets, and more—to deliver accurate, context-aware responses. This innovative integration enables the system to maintain relevant context across user interactions and improve answer precision.
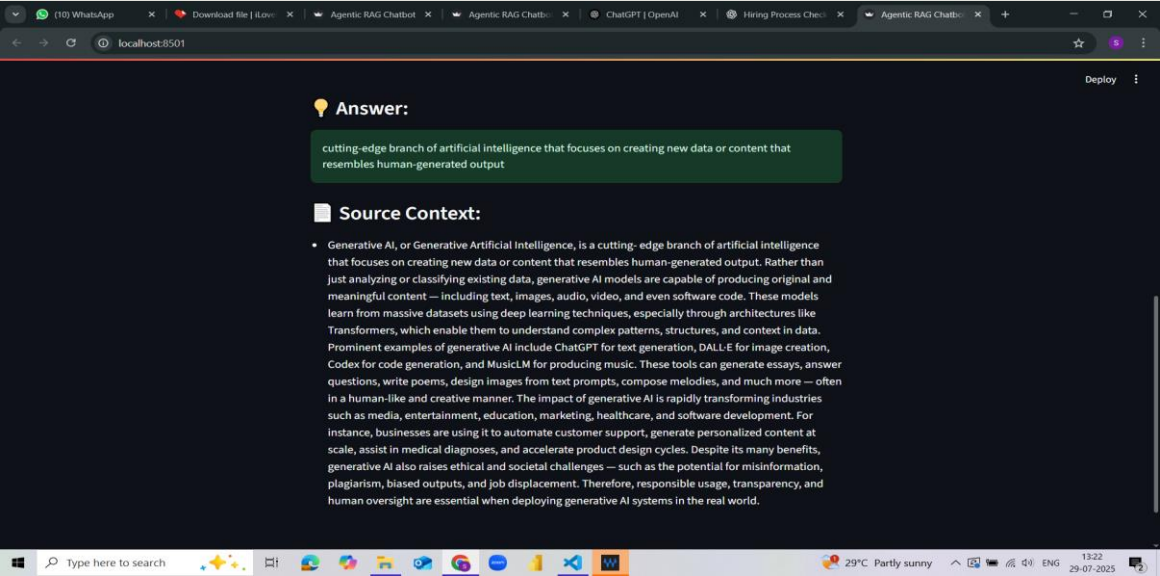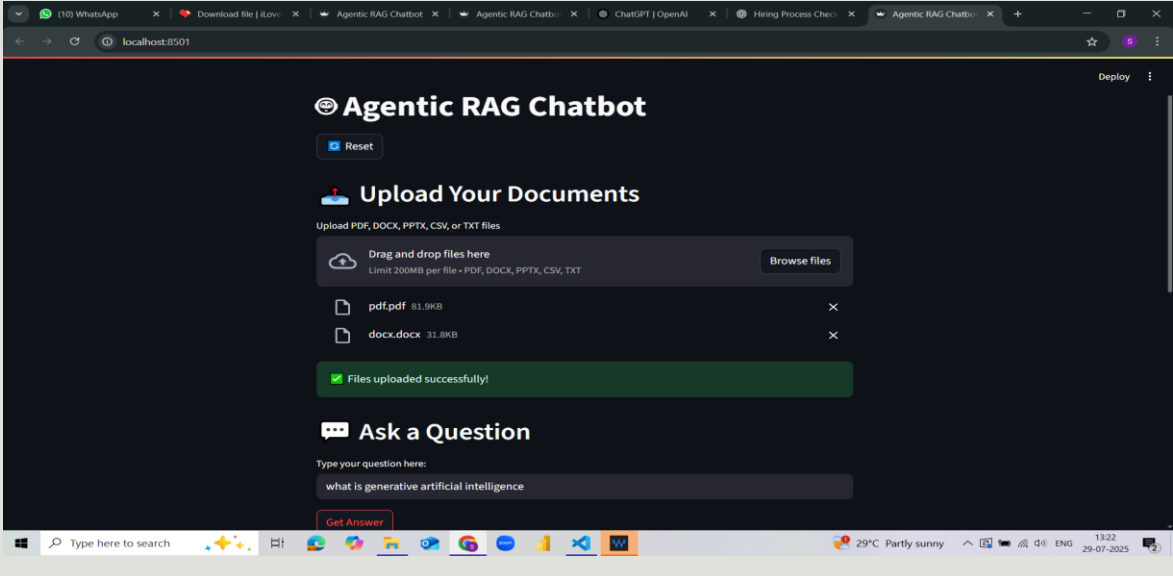
## MCP protocol

```python
import uuid

def create_mcp_message(sender, receiver, msg_type, payload, trace_id=None):
    return {
        "sender": sender,
        "receiver": receiver,
        "type": msg_type,
        "trace_id": trace_id or str(uuid.uuid4()),
        "payload": payload
    }
```

| Category | Technology | Purpose |
|----------|-----------|---------|
| Frontend (UI) | Streamlit | Chat interface, file uploads, user interaction |
| Backend Logic | Python | Core development language |
| Agents | Custom Python Classes | IngestionAgent, RetrievalAgent, LLMResponseAgent |
| Messaging Layer | Model Context Protocol (MCP) | Message-passing between agents (in-memory / REST) |
| Document Parsing | PyMuPDF, python-docx, python-pptx, pandas | Parsing PDF, DOCX, PPTX, CSV, TXT files |
| Embeddings | Hugging Face Transformers | Generating vector representations of chunks |
| Vector Store | Vector DB with FAISS | Storing and retrieving document embeddings |
| LLM | google/flan-t5-base | Answer generation using retrieved context |
| File Handling | Python OS / tempfile / base64 | Upload, read, and manage user documents |
| Deployment | Streamlit | App deployment and demo hosting |
| Version Control | Git + GitHub | Code management and collaboration |

# Results:

# Challenges Faced

**Repetitive Content in Source Documents**

Several documents contained repeated paragraphs or duplicated content across different formats (e.g., same text in PDF and PPTX).

This led to redundant chunks being passed to the LLM, increasing token usage and sometimes causing repetitive or verbose responses.

**Solution:** Implemented a deduplication step during document parsing to filter out repeated text before embedding and retrieval.