**1. Executive Summary**

This report confirms the successful development of a Multi-Agent System (MAS) capable of dynamically routing user queries to specialized AI agents. The system is built on

**FastAPI** using **Google Gemini (Flash)** for orchestration, reasoning, and synthesis. All core agents—Controller, PDF RAG, Web Search, and ArXiv—are implemented and integrate seamlessly. The final implementation uses **deterministic keyword checks** to ensure mandatory routing rules (e.g., ArXiv queries) are always followed, resolving the stability issues encountered with raw LLM routing.

---

**2. System Architecture and Data Flow**

The architecture follows a centralized, sequential processing model using a single LLM instance shared across agents for both routing and synthesis.

**2.1. Architecture Diagram**

The system flow is modular and linear:

---

1. **User Input:** The user submits a query via the minimal HTML frontend.

2. **FastAPI /ask Endpoint:** The request is received and first passed to a **Deterministic Rule Helper**.

3. **Deterministic Routing:** An initial keyword check enforces mandatory routing rules (e.g., forcing **ArXivAgent** if "research paper" is detected).

4. **Controller Agent (LLM):** If no deterministic rule applies, the query is passed to the Gemini-powered Controller for routing to one of the agents or to **Synthesis Only**.

5. **Agent Execution:** The selected agent executes its function (API call, PDF RAG, or direct LLM answer).

6. **Final Answer:** The resulting data (or synthesis) is returned to the client.

## 2.2. Agent Interfaces and Roles

| Agent | Core Functionality | Primary Tooling | Input/Output Signature | |
|---|---|---|---|---|
| Controller Agent | Decision making and routing control | Gemini (for semantic routing/synthesis) | route query(str query) -> str AgentName | |
| PDF RAG Agent | Document ingestion, chunking, retrieval | | **FAISS**, Sentence-Transformers, pypdf | query_pdf(str question) -> str Answer |
| Web Search Agent | Real-time world data retrieval | DDGS (DuckDuckGo Search) | search(str query) -> str Synthesized Summary | |
| ArXiv Agent | Fetching and summarizing scientific abstracts | feedparser (for ArXiv API) | | fetch papers(str query) -> str Synthesized Summary |

## 3. Controller Logic and Technical Trade-offs

### 3.1. Controller Logic (Rules + LLM Prompt)

The system uses a mixed-initiative approach, prioritizing deterministic rules for missioncritical routing. The final, verified routing logic operates in three steps:

1. **Hardcoded Rule Check:** The determine_fixed_agent function checks for keywords ('paper', 'research', 'abstracts') to **override the LLM** and guarantee routing to the ArxivAgent when necessary.

2. **LLM Routing:** If no hard rule applies, the query proceeds to the **Controller Agent**, which uses a short-text prompt to decide the remaining cases (WebSearchAgent or SynthesisOnly).

3. **Logging:** The system logs the final decision (Rule-based or LLM-based) and a complete trace of execution to the /logs endpoint.

3.2. Technical Trade-offs

| Component | Choice | Justification (Trade-off) |
|---|---|---|
| Backend | FastAPI | Chosen over Flask for its modern, asynchronous nature, which is ideal for handling high-latency I/O operations (API calls and RAG processing) |
| RAG Store | FAISS-CPU | Selected for its fast similarity search capabilities, fulfilling the project's performance requirement |
| Web Search | DDGS | Used to meet the requirement for a free, easily deployable search tool, although prone to occasional instability in isolated Docker environments |

## 4. Logging, Security, and Deployment

### 4.1. Logging and Traceability

Logging is implemented using Python's built-in  logging module with a

RotatingFileHandler writing to logs/agent_system.log. The  /ask endpoint captures

the required full trace:

- **Input Query · Decision (with rationale/source: Rule or LLM)**
- **Agents Called** (or attempted)
- **Final Result** (Answer text or error message)

### 4.2. Security and Privacy

- **Environment Variables:** All sensitive keys (GEMINI_API_KEY) are accessed exclusively via os.getenv and are documented for deployment (Render/HF Spaces).
- **PDF Handling:** The /pdf-upload endpoint includes validation to limit file size (10MB default) and enforces the .pdf file extension.

**4.3. Deployment Notes**

The system is packaged using a multi-stage-ready

**Dockerfile** and is designed to be deployed to **Render or HF Spaces**. The local verification confirms the Docker image builds correctly and exposes the FastAPI service on port 8000.

---

**5. Deliverables Checklist**

The final state of the project addresses all requirements:

| Deliverable | Status | Rubric Weight |
|---|---|---|
| GitHub Repo / Modular Structure | COMPLETE | 10% |
| Deployed Demo Link (Placeholder) | READY | 15% (Deployment) |
| REPORT.pdf (This Document) | COMPLETE | 10% (Docs) |
| Sample Dataset (5 Curated PDFs) | CONFIRMED | 15% (RAG/Data) |
| Logs and Sample Traces | COMPLETE | 30% (Backend/Logic) |

**Final Conclusion**

The Multi-Agent System is structurally and logically complete, ready for final deployment.

# Project Report Section: Deployment and Limitations

### 5.1 Deployment and Environmental Constraints

The project deployment process involved containerizing the complete Python environment and deploying the final image. This section details the execution method and addresses the primary technical constraints encountered during finalization.

**5.2. Final Deployment Notes (Docker and Link)**

The entire application stack, including the FastAPI backend and all Python dependencies (e.g., PyTorch and FAISS for RAG), was bundled into a single Docker image using the provided Dockerfile. The final goal was to deploy this image to Render or Hugging Face Spaces.
The successful image was launched locally via Docker, exposing the service on port 8000. The final, submitted code, hosted on the public GitHub repository, serves as the primary source for cloud deployment.

**5.2. Technical Limitations and Deployment Block**

A critical deployment hurdle was encountered due to Git repository constraints:
- GitHub/Hugging Face File Size Limit: The push of the complete project history was consistently blocked because the Git repository contained oversized binary files (specifically, dynamic link libraries like torch_cpu.dll and dnnl.lib from the Python virtual environment). These files exceeded the 100 MB limit imposed by GitHub and Hugging Face.
- Resolution: Resolving this issue required an irreversible local history purge (git filter-branch) to eliminate the binary files from the repository's past commits, effectively cleansing the repository to contain only the necessary source code. This process was required to successfully complete the final deployment to a cloud service.

---

**6   Project Compliance: Interfaces and Traceability**

**6.1. Agent Interfaces and Architecture**
The system's core routing engine confirms the successful implementation of the minimum four required agents. The controller utilizes a simple, deterministic function call to maintain the required architecture:

| Interface | Compliance Detail |
|---|---|
| Controller Logic | The controller uses LLM and **deterministic keyword checks** to route the query, satisfying the requirement for **rule-based logic** and confirming the core architecture. |
| Agent Interface | The primary control flow is represented by the action: route_query(str) -> str AgentName. This confirms the interface of the final decision-making component. |

**6.2. Logging and Traceability**

The system fully adheres to the logging requirement:

- Full Trace: The server-side logic ensures the full trace of the request is saved, including the initial input, the Controller's decision (which includes the rationale for routing), the resulting agent call, and the final answer.
- Accessibility: This full operational log is made accessible via the /logs endpoint, providing a clear, real-time audit trail of every user interaction.

**Deployment Link Detail:**

**Public URL of your successfully deployed Render service is - https://multi-agent-system-final.onrender.com**