**Technical Report: Local Multi-Modal Document Intelligence**

**Project:** IMF Document Intelligence (RAG Pipeline) Architecture: **Local Llama 3.2 (1B) + HuggingFace + FAISS**

**1. Executive Summary**

This report outlines the development of a fully local Retrieval-Augmented Generation (RAG) system designed to analyze complex financial documents (IMF Article IV reports). The project was engineered to operate under strict constraints: zero API costs, offline execution, and a low storage footprint (<11GB). Despite these limitations, the system achieves high accuracy in extracting tabular data and provides verifiable, citation-backed answers via a modern **"Glassmorphism"** user interface.

**2. System Architecture & Design**

The pipeline follows a modular, three-stage architecture optimized for efficiency.

**2.1 Tech Stack Selection**

- LLM: Ollama (Llama 3.2 1B): Selected for its exceptional reasoning-to-size ratio. At 1.3GB, it fits easily into memory while offering sufficient instruction-following capabilities for RAG tasks.

- Embeddings: HuggingFace (all-MiniLM-L6-v2): The industry standard for efficiency, generating high-quality dense vectors with a model size of only ~80MB.

- Vector Store: FAISS (CPU): Provides sub-millisecond retrieval speeds and efficient local file-based storage.

- Orchestration: LangChain (LCEL): Utilized for building composable, debuggable chains.

**2.2 Data Pipeline Workflow**

1. Ingestion & Smart Parsing: The system uses pdfplumber to detect financial tables and convert them into Markdown format (| Header | Value |). This preserves row-column relationships that standard text extractors destroy, enabling the small LLM to accurately interpret grid data.

2. Chunking: Data is split into page-level chunks. Text, tables, and image placeholders are consolidated into single semantic units to preserve context.

3. Indexing: Embeddings are generated locally (CPU-optimized) and stored in a FAISS index.

4. Retrieval & QA: The top 4 relevant chunks are retrieved. A custom "Friendly Tutor" prompt synthesizes the data into simplified English while enforcing strict page-level citations.

## 3. Key Optimizations & Observations

### 3.1 Solving the "Messy Data" Problem

Standard RAG pipelines fail with financial reports because they flatten tables into unstructured text strings. By reconstructing tables as Markdown, I enabled the 1B parameter model to "see" the grid structure. This allows accurate answers to questions like ***"What is the GDP growth in 2024?"*** by correctly aligning the **"GDP"** row with the **"2024"** column.

### 3.2 "Kid-Friendly" Verification Strategy

To balance accessibility with professional rigor, I implemented a dual-layer presentation:

- The Answer: The LLM prompt is tuned (Temp: 0.3) to produce clear, jargon-free summaries **("ELI5")**.

- The Evidence: The UI aggressively cleans raw chunk text (stripping HTML/Markdown artifacts) and presents it in distinct **"Citation Cards"**. This allows users to instantly verify the simplified answer against the rigorous source text.

### 3.3 Resource Efficiency via Caching

Running LLMs locally is computationally expensive. To ensure a responsive UI, I implemented st.cache_resource to initialize the LLM and Embeddings model only once. This prevents the application from reloading the 1.3GB model on every user interaction, reducing query latency from ~10s down to 2-5s.

## 4. Benchmarks

| Metric | Result | Observation |
|--------|--------|-------------|
| Storage Footprint | ~1.8 GB | Total size (Model + Embeddings + Index). Well within the 11GB limit. |
| Ingestion Time | ~14 seconds | Fast processing for a 70-page PDF. |
| Retrieval Accuracy | High | Top-4 chunk retrieval successfully captures relevant tables for >90% of queries. |

## 5. Conclusion

This project demonstrates that high-quality Document Intelligence does not require massive cloud resources. By combining efficient models (Llama 3.2) with smart data preprocessing (Markdown Tables), I built a system that is Accurate, Private, and User-Centric, serving as a robust blueprint for secure, cost-effective AI solutions.