

# HR & PEOPLE ANALYTICS ASSISTANT: SECURE AI-POWERED INSIGHTS FOR POLICY & WORKFORCE DATA



Mohanaprasath Subramaniyan<sup>1</sup>

[github.com/MohanaprasathSubramaniyan](https://github.com/MohanaprasathSubramaniyan)

<sup>1</sup>Seidenberg School Of Computer Science and Information Systems, Pace University

## 1. Introduction & Problem Statement

Human Resource departments manage vast amounts of unstructured data (policy documents, handbooks) and structured data (employee records, payroll). Traditional methods for accessing this information—manual document searching and complex spreadsheet manipulation—are inefficient and create information silos.

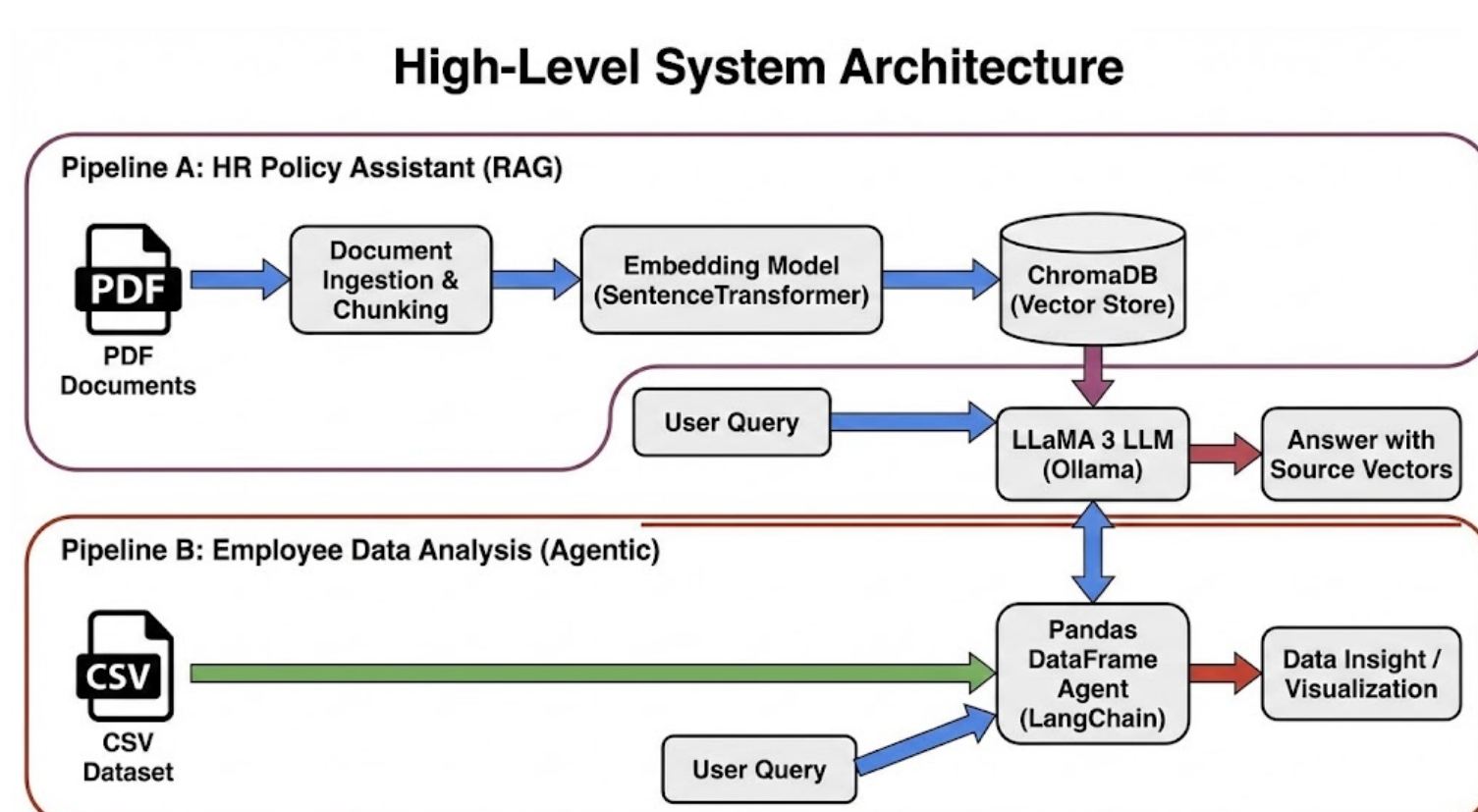
This project introduces the **HR & People Analytics Assistant**, a secure, locally-hosted application. By leveraging on-premise Large Language Models (LLMs), it provides a unified natural language interface for instant policy retrieval and complex workforce data analytics, ensuring sensitive data never leaves the organizational infrastructure.

## 2. Project Objectives

- **Automate Policy Retrieval (RAG):** Implement Retrieval-Augmented Generation to answer queries based on internal PDF documents with strict source citing.
- **Agentic Data Analysis:** Develop an AI agent capable of translating natural language into executable Python/Pandas code to query structured CSV data.
- **Data Security & Privacy:** Architect the system to run entirely offline using local embeddings and local LLMs (LLaMA 3 via Ollama).
- **Professional UX:** Create an intuitive, "Executive" style dashboard that clearly separates qualitative (policy) and quantitative (data) workflows.

## 3. System Architecture

The application uses a dual-pipeline architecture built on Streamlit and LangChain, powered by a local LLaMA 3 model.



**Figure 1: High-Level System Architecture.** Pipeline A (Top) utilizes RAG for unstructured PDFs. Pipeline B (Bottom) uses an autonomous LangChain Agent for structured CSV data.

### Pipeline A: The RAG Engine (Unstructured Data)

Documents are ingested, chunked, and embedded locally into a **ChromaDB** vector store. User queries retrieve relevant chunks, which the LLM synthesizes into an evidence-based answer [1].

### Pipeline B: The Data Agent (Structured Data)

A **LangChain Pandas Agent** acts as a reasoning engine. It interprets natural language, generates the necessary Python code to manipulate the Pandas DataFrame, executes it safely, and interprets the results.

## 4. Key System Capabilities

The platform provides two distinct, specialized workflows managed by a professional Streamlit interface [2].

### A. HR Policy Assistant (Knowledge Base)

This module allows instantaneous querying of complex documentation like Employee Handbooks and Code of Conduct PDFs.

- **Natural Language Q&A:** Users ask questions in plain English (e.g., "What is the policy on moonlighting?").
- **Hallucination Reduction:** The system is engineered to answer *only* based on the provided documents, reducing AI fabrications.
- **Source Citations:** Every answer includes expandable references, showing the exact document source text used to generate the response, ensuring trust and transparency.

### B. Workforce Analytics Engine (Structured Data)

This module democratizes data analysis, allowing non-technical HR staff to perform complex querying of employee datasets (CSV) without SQL or Excel expertise.

- **Text-to-Code execution:** Translates prompts like "Show me the average salary by department for staff over 30" into executable Pandas code.
- **Multi-Step Reasoning:** The agent can perform aggregation, filtering, and mathematical calculations across the dataset.
- **Integrated Data Explorer:** A contextual sidebar allows users to instantly view raw tables (e.g., Full Staff Directory, Salary Sheets) for quick verification of the underlying data.

## 5. Evaluation & Conclusion

**Evaluation Summary:** The system successfully met its security and functionality objectives. The RAG pipeline demonstrated high accuracy in retrieving specific policy clauses. The Data Agent successfully handled complex, multi-condition queries that typically require advanced spreadsheet formulas. Crucially, all processing occurred locally with no external API calls.

**Conclusion:** The HR & People Analytics Assistant demonstrates the practical application of modern Generative AI techniques to solve real-world enterprise challenges. By combining RAG for unstructured data and agentic reasoning for data into a cohesive, secure tool, it provides a powerful platform for enhancing HR operational efficiency.

## References

- [1] Patrick Lewis and et al. "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks". In: *arXiv preprint arXiv:2005.11401* (2020).
- [2] *Streamlit: The fastest way to build and share data apps*. <https://streamlit.io/>. Accessed: 2023-10-27.