# Llama Cloud -
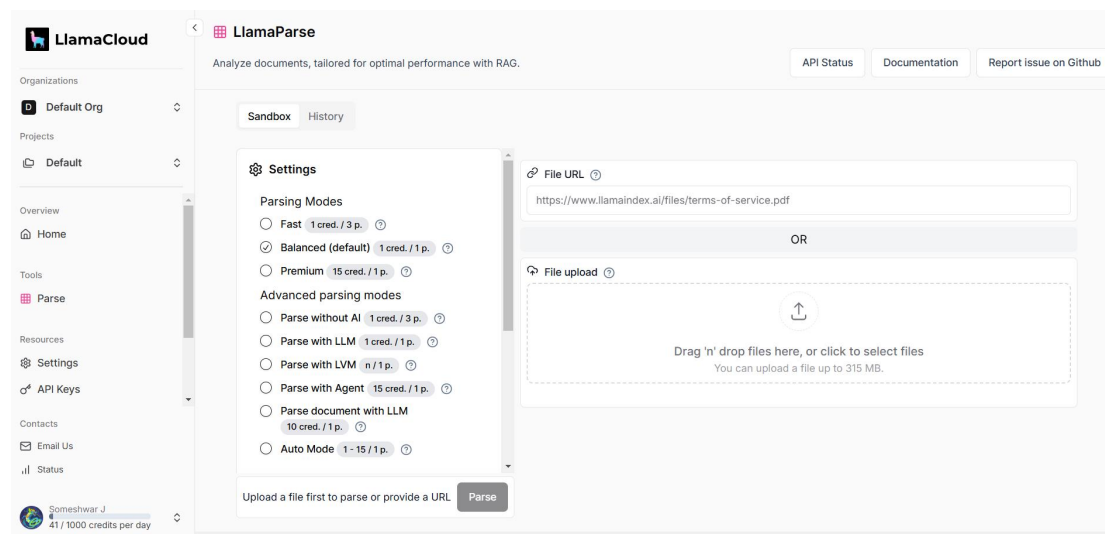# RAG implementation using Llama Parsing with Trulens Evaluation

This project implements a simple Retrieval-Augmented Generation (RAG) pipeline. It involves parsing an uploaded document, setting up a RAG system, and evaluating its performance using TruLens.

## 1. Parsing the Uploaded Document

The first step in the pipeline is document parsing. The uploaded document is processed to extract text content, which is then used in the RAG pipeline. The parsing mechanism ensures the extracted content is clean and structured for retrieval.
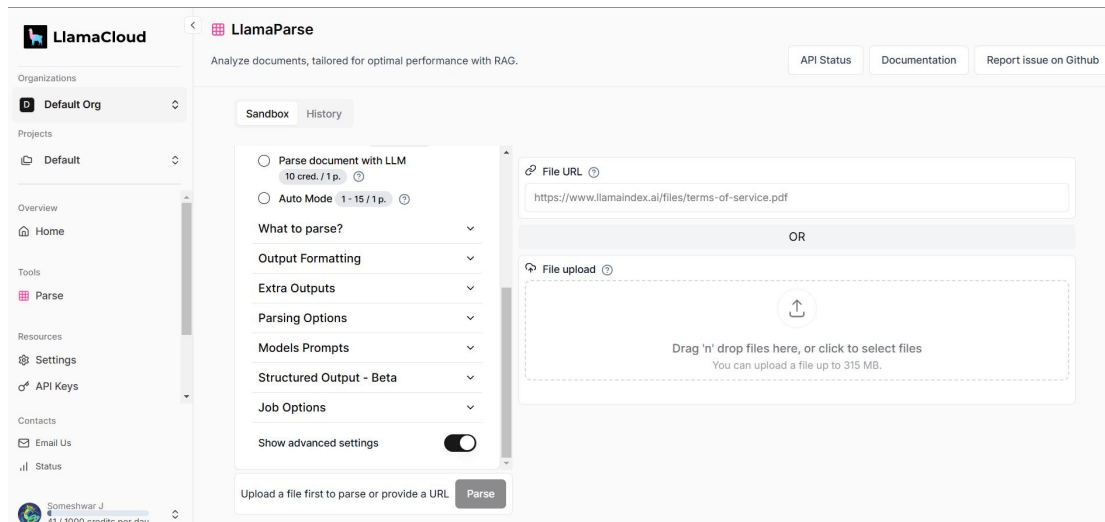
**Steps Involved:**

- Uploading the Document - The user provides a document (e.g., PDF, TXT).
- Extracting Text - The document is parsed to extract textual content for further processing.
- Preprocessing - Cleaning and structuring the text to ensure quality input for retrieval.



## Llama Cloud - Document Parsing and Agent Builder using Ollama

- Can parse, extract, and index a million of your most complex PDFs, Powerpoints, and more.
- Build Agent workflows on top of this.
- Link to create explore: [Llama Cloud](#)
- Documentation Link: [LlamaParse Documentation](#)

LlamaParse with Advanced Options

## 2. Building the RAG Application

Retrieval-Augmented Generation enhances language model responses by incorporating relevant retrieved documents. The RAG system consists of the following components:

**Components:**

- Vector Database - Stores document embeddings for efficient retrieval.
- Embedding Model - Converts text into numerical representations.
- Retriever - Fetches the most relevant text passages based on queries.
- LLM (Large Language Model) - Generates responses using retrieved information.

**Workflow:**

- Indexing the Document - The extracted text is converted into embeddings and stored in a vector database.
- Local Storage for Indexing - Since LlamaIndex is not free, the system uses local storage to save vectors and embeddings instead of a cloud-based solution.
- Retrieval - When a query is made, the retriever fetches relevant text chunks.
- Querying with LlamaIndex - The retriever utilizes the query capabilities from the LlamaIndex Python package to find the most relevant information.
- Generation - The retrieved text is passed to the LLM, which generates a final response.
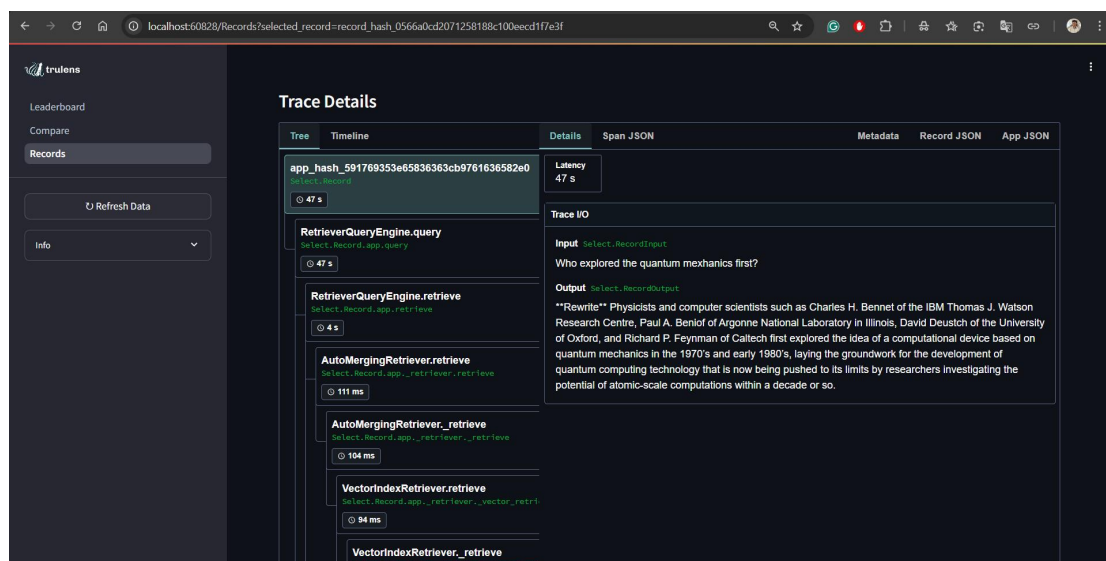
## 3. Evaluating with TruLens

TruLens is used to assess the quality and reliability of the RAG system. It helps measure response relevance, faithfulness to retrieved data, and overall effectiveness.

**Evaluation Metrics:**

- Relevance - How well the retrieved text matches the query.
- Faithfulness - Whether the generated response accurately reflects the retrieved information.
- Latency - Measuring system response time.

**Steps:**

- Logging RAG Execution - Tracking retrieval and generation steps.
- Applying TruLens Metrics - Using TruLens to evaluate output quality.
- Interpreting Results - Analyzing performance and identifying areas for improvement.



Trulens UI

✧ Click here to explore the [notebook](notebook)