# Design and Implementation of an AI-Integrated Web Application Using Large Language Models and Retrieval-Augmented Generation

Samuel Rivera

*Abstract*—This report presents the design and implementation of an AI-integrated web application that combines client-side machine learning, server-side large language models (LLMs), and Retrieval-Augmented Generation (RAG). The system demonstrates multiple artificial intelligence paradigms, including image classification using TensorFlow.js, natural language processing with OpenAI models, and document-based reasoning using LangChain. A web-based user interface enables users to upload PDF resumes, which are then analyzed by a RAG pipeline to generate both praising and critical evaluations.

*Index Terms*—Artificial Intelligence, Large Language Models, LangChain, Retrieval-Augmented Generation, TensorFlow.js, Web Applications.

## I. Introduction

Artificial intelligence is increasingly embedded in modern software systems, enabling applications to perform complex tasks such as perception, reasoning, and decision support. This project demonstrates a full-stack AI web application that integrates client-side inference, server-side LLMs, and document-grounded reasoning using Retrieval-Augmented Generation (RAG).

## II. Client-Side AI with TensorFlow.js

TensorFlow.js enables machine learning models to run directly in the browser. In this project, a pre-trained MobileNet model was used to classify images without requiring a backend server.



Fig. 1: Client-side image classification using TensorFlow.js and MobileNet.

The application was extended to allow users to upload new images and trigger classification dynamically.

## III. Large Language Models with OpenAI

The backend integrates OpenAI's JavaScript SDK to interact with large language models. Multiple prompts and model configurations were tested to compare response quality and behavior.
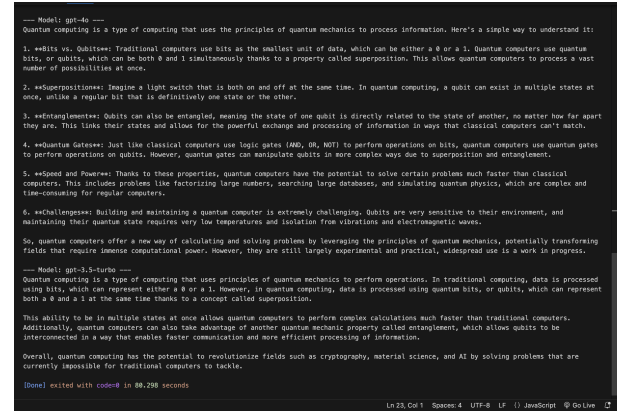


Fig. 2: Terminal output demonstrating OpenAI LLM responses to different prompts.

This component demonstrates how LLMs can be used for content generation and conversational tasks.

## IV. LangChain Guardian Architecture

A LangChain-based guardian architecture was implemented to explore AI safety and control. The system includes a Red Team model, a Worker model, and a Guardian model that evaluates prompts and responses for safety.

This layered approach illustrates how AI systems can enforce safety constraints.

## V. Retrieval-Augmented Generation

Retrieval-Augmented Generation enhances LLM outputs by grounding responses in external documents. Uploaded PDF resumes are parsed, chunked, embedded, and stored in a vector database. Relevant chunks are retrieved and injected into prompts.

Two prompts generate complementary outputs: a praising analysis and a critical analysis.
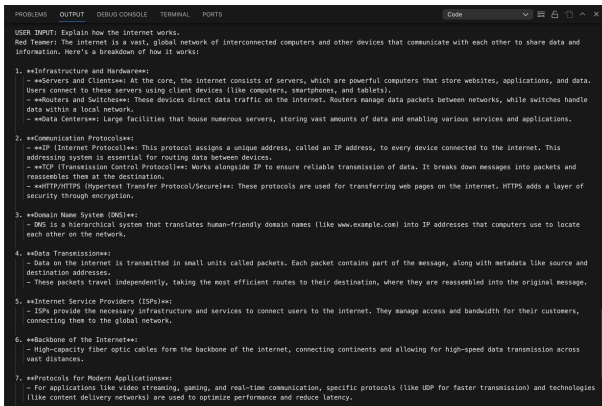
Augmented Generation, the system highlights practical AI deployment patterns and real-world engineering challenges.

REFERENCES

[1] TensorFlow.js Documentation. https://www.tensorflow.org/js
[2] OpenAI API Documentation. https://platform.openai.com/docs
[3] LangChain Documentation. https://js.langchain.com
[4] Express.js Documentation. https://expressjs.com
[5] A. G. Howard et al., "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications," arXiv:1704.04861, 2017.

Fig. 3: LangChain Guardian workflow showing Red Team, Worker, and Guardian decisions.



Fig. 4: Retrieval-Augmented Generation pipeline for resume analysis.

## VI. WEB-BASED USER INTERFACE

A web-based graphical user interface allows users to upload PDF resumes and view AI-generated analyses. The backend is implemented using Node.js and Express, with Multer handling file uploads.
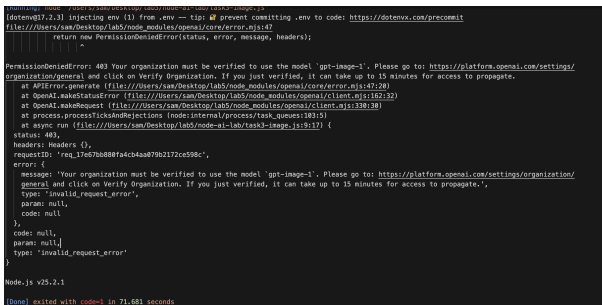


Fig. 5: Web interface for uploading a resume PDF.

## VII. CHALLENGES AND RESULTS

Several challenges were encountered, including dependency conflicts, ES module compatibility issues, and operating system file metadata problems. Despite these challenges, the final system successfully integrated all required components. The RAG-based analysis produced more relevant and grounded outputs compared to standalone LLM responses.

## VIII. CONCLUSION

This project demonstrates the integration of modern AI technologies into a full-stack web application. By combining TensorFlow.js, OpenAI LLMs, LangChain, and Retrieval-