# Final Project Report – Suman Sahu

# Financial Report Generation

## 1. Project Overview

This project leverages an Agent built using the LlamaIndex framework, HuggingFace embeddings, and the Groq-hosted llama3-70b-8192, language model. The goal is to extract and compare top-level financial data (assets and liabilities) of Apple and Tesla from structured reports using natural language queries. By automating this analysis with AI tools, we aim to demonstrate how Large Language Models (LLMs) can assist in financial reasoning and business intelligence tasks.

## 2. Objectives

- Set up an AI agent pipeline for document understanding using LlamaIndex and Groq LLM.
- Extract top-level assets and liabilities for:
    - Tesla in 2019-2023
    - Apple in 2019-2023
- Create a VectorIndex for context retrieval using both companies data.
- Setup the generation pipeline.
- Create Agent tool pipeline using LLM which supports function calling.
- Summarize findings and draw insights using AI-driven outputs.

## 3. Setup & Initialization

The project begins with the initialization of essential components:

- `LlamaCloudIndex`: for managing and querying documents.
- Environment setup using `os`, `nest_asyncio`, and `transformers`.
- Groq-hosted **llama3-70b-8192** for generating structured and coherent financial summaries.
- Input data: Text-based financial documents stored under `/data`.

# 4. Tools & Technologies

| Component | Description |
|---|---|
| **LlamaIndex** | Core framework for indexing, retrieval, and agent orchestration using LLMs. |
| **LlamaCloud** | Managed vector store enabling efficient file-level and chunk-level retrieval. |
| **Groq LLM API** | Ultra-fast inference backend powering **LLaMA3-70B-8192** for both generation and agent pipelines. |
| **Hugging Face Embeddings** | Used to convert document content into vector embeddings for semantic search. |
| **Python Libraries** | Includes `asyncio`, `transformers`, `os`, and `pydantic` for environment setup and model integration. |

# 5. Project Structure

## 5.1 Setup & Initialization

- Importing libraries (`llama_index`, `transformers`, `groq`, etc.)
- Loading tokenizer and LLM
- Setting configurations and embeddings

## 5.2 Document Loading

- Using `SimpleDirectoryReader` to ingest PDF or text files
- Preprocessing chunks for indexing

## 5.3 Index Creation

- Building `VectorStoreIndex` from loaded documents
- Creating retrievers at chunk and file levels

## 5.4 LLM setup

- Initialized the **LLaMA3-70B-8192** model via the **Groq API** for high-speed LLM inference.
- Configured model parameters like **temperature** and **API key** for controlled generation.
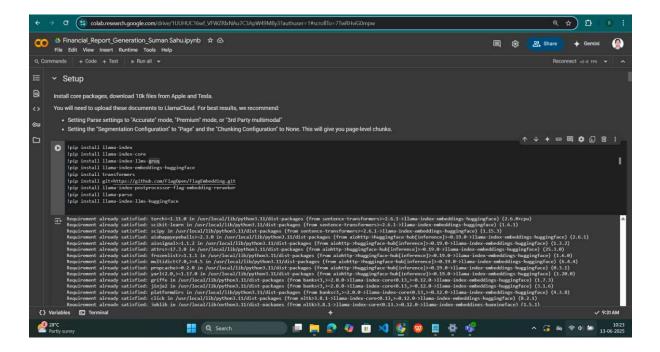
## 5.5 Agent setup

- Initialized agent so it will intelligently select between document-level and chunk-level retrieval based on the query.
- To generate a structured report by combining relevant text and tables for clear financial insights.
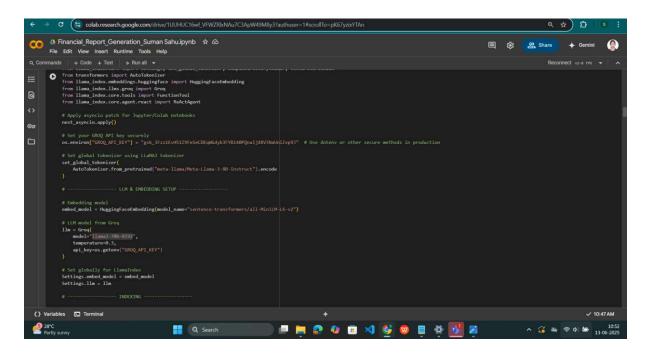
## 5.6 Query Handling

- Asking complex questions like:
  - "Compare Apple and Tesla financials 2020–2023"
- Receiving structured LLM responses
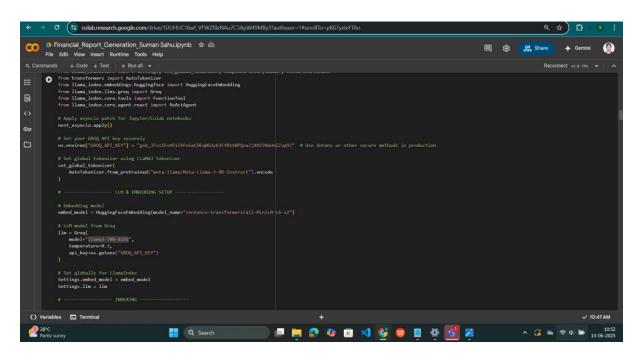
# 6. Project Setup

## 6.1 Environment Initialization
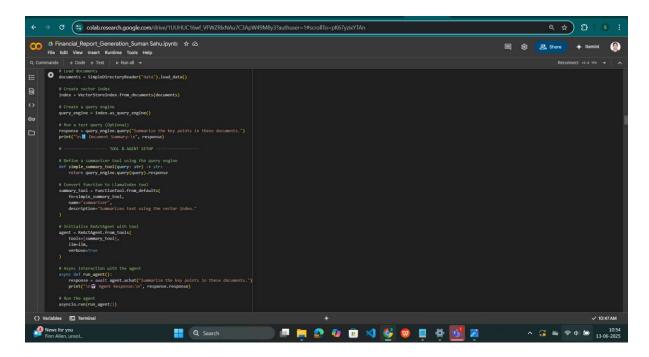
## 6.2 Tokenizer and LLM setup



## 6.3 Connecting Groq LLM
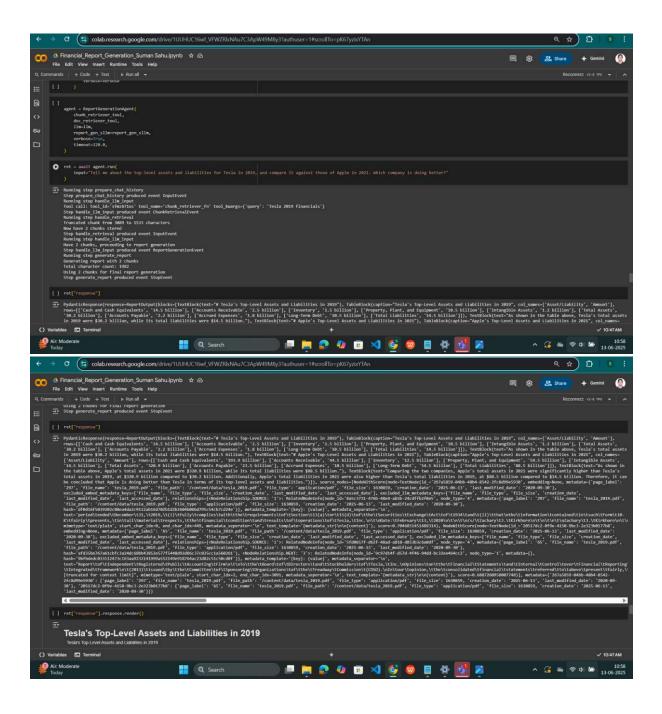
# 7. Vector Index creation

We placed structured financial reports for Apple and Tesla inside a folder named `data`.
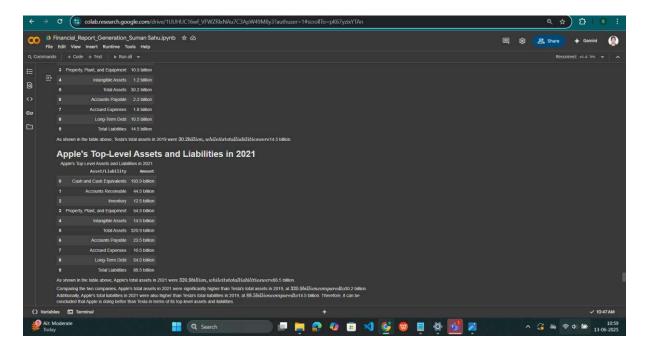
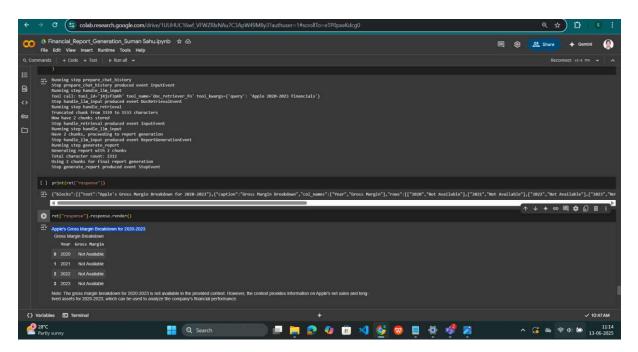# 8. Agent Initialization and Query Execution

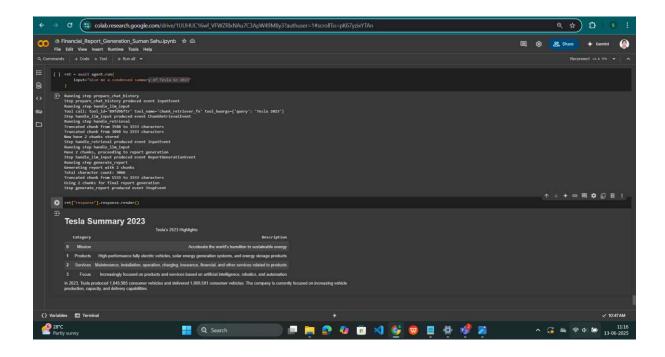We asked the agent three questions:

1. **Comparative Financial Analysis**

# 9. Report Generation



## 2. Apple's Gross Margin Breakdown for 2020-2023



## 3. Tesla Summary 2023

# 10. Conclusion

This project showcases the practical strength of custom AI agents in automating advanced financial analysis. By building a **custom `ReportGenerationAgent`** that integrates **LlamaIndex** for retrieval, **Hugging Face embeddings** for semantic indexing, and **Groq-hosted LLaMA3-70B-8192** for both structured generation and tool-augmented reasoning, we demonstrated the ability to:

- Ingest and process structured financial reports effectively
- Retrieve relevant content at chunk and document levels using modular tools
- Generate accurate, structured insights combining textual analysis and tabular reasoning
- Compare financial performance across companies and years with contextual depth

This project confirms the capability of **function-calling, structured-output agents** to serve as intelligent financial assistants—automating repetitive yet complex tasks with both speed and precision, making them invaluable in decision-critical environments.