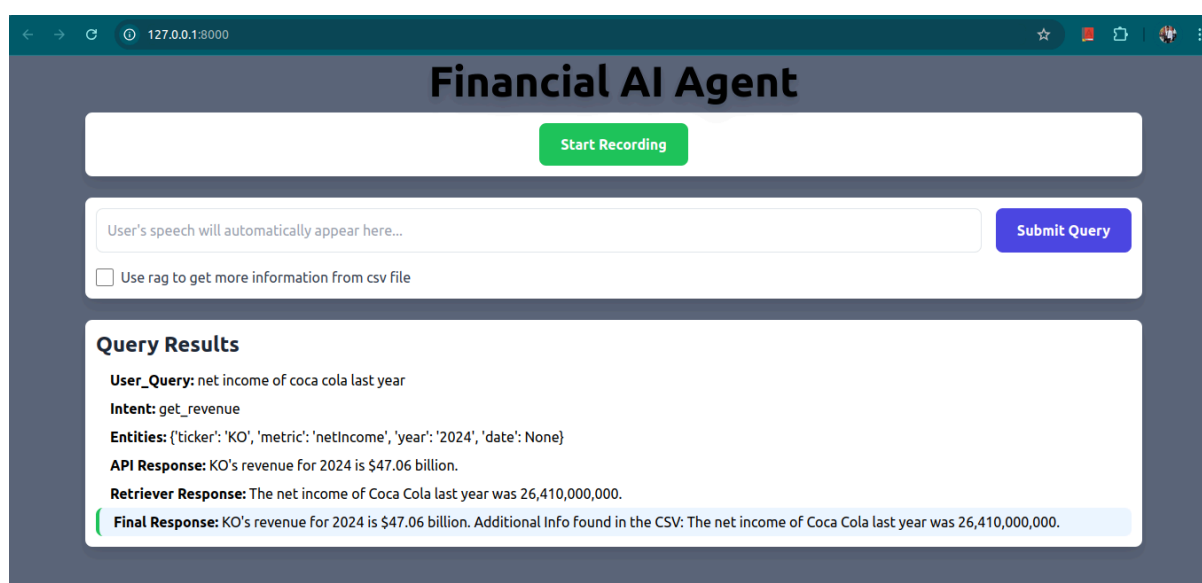


# Documentation of Financial AI Agent

## Overview

The Financial AI Agent is an intelligent system designed to query and retrieve critical financial data from bank and company financial statements using voice recognition. It enables users to access key financial metrics (e.g., net income, revenue) and manage company assets efficiently. The system integrates voice input, natural language processing (NLP), data retrieval from CSV files and the web, and response generation to provide accurate and actionable financial insights. **The project is fully functional on a CPU machine.**



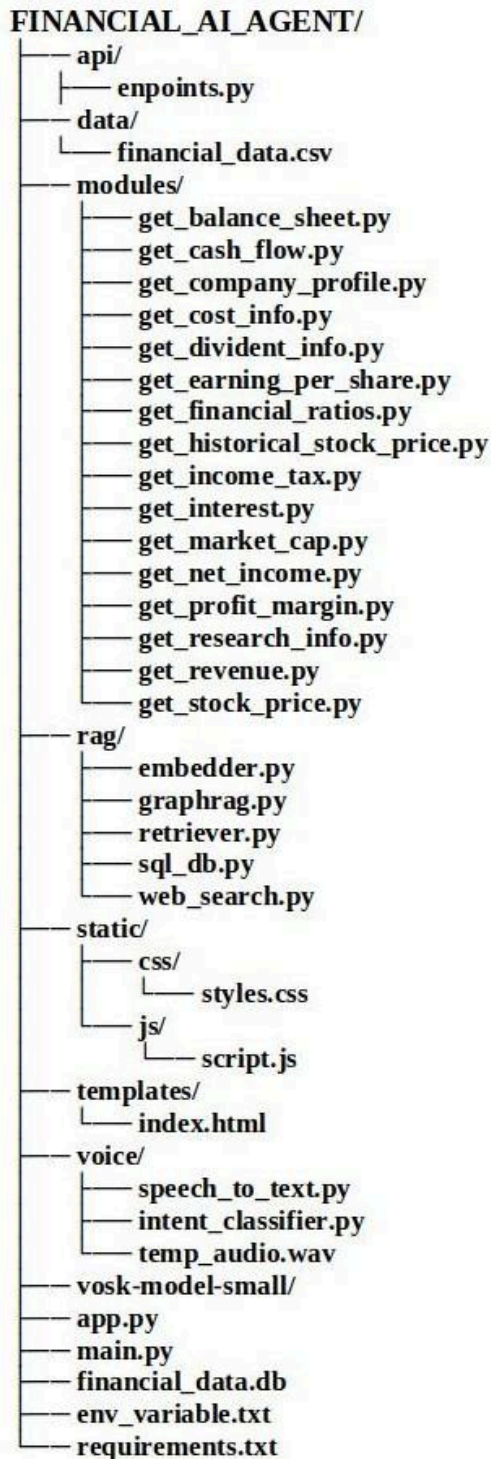
## Key Features

- **Voice-Activated Financial Queries:** Enables users to query financial data (e.g., net income, revenue) using voice commands, leveraging offline speech recognition with Vosk for privacy and reliability.
- **Multi-Source Data Retrieval:** Seamlessly integrates multiple data sources—Financial Modeling Prep API, local CSV files (financial\_data.csv), and web searches via DuckDuckGo—to ensure high availability of financial information.
- **Dynamic Financial Metric Detection:** Automatically identifies financial metrics (e.g., "net income", "revenue") in user queries and tailors searches to retrieve relevant data, supporting a wide range of financial inquiries.

- **RAG Integration with FAISS:** Uses FAISS for efficient similarity search and retrieval of financial documents, enhancing the accuracy and relevance of responses through vector-based search.
- **SQLite Database Integration:** Stores and manages financial data in an SQLite database, enabling efficient querying and persistence of financial metrics and user data.
- **Intelligent Name Entity Recognition:** Maps company names and other metrics to specific symbols using a predefined dictionary or mapping which enhances accuracy by searching financial\_data.csv for matches with  $\geq 80\%$  similarity, ensuring robust entity extraction.
- **Privacy-Focused Web Search:** Utilizes duckduckgo\_search for web queries, a privacy-focused solution that requires no API key or payment details, making it accessible and secure.
- **Numerical Data Prioritization:** Prioritizes web search results containing numerical values (e.g., "\$10.631B") to deliver precise and actionable financial insights to the user.
- **Modular and Extensible Design:** Organized into modular components (speech-to-text, intent classifier, retriever, etc.), allowing easy integration of new data sources or features (e.g., additional financial APIs, global datasets).
- **FastAPI-Powered User Interface:** Provides a modern, asynchronous web interface using FastAPI, allowing users to interact with the financial agent via a browser, view financial data, and visualize reports.

# Project Architecture

The Financial AI Agent is structured as a modular system with distinct components, each handling a specific functionality. The architecture follows a pipeline approach: voice input is processed, intent and entities are extracted, data is retrieved from api calling or local directories or web search, and a response is generated for the user.



# Project Mechanism Overview

This project provides an interactive user interface for the **Financial AI Agent**, developed using **FastAPI** and **Uvicorn**. The API allows users to manage the entire codebase and interact with the system for voice-activated financial data retrieval. Below is a detailed explanation of how the system processes user queries:

## 1. Voice Input & Transcription

- The user initiates voice input by clicking the **Start Recording** button. Recording continues until the **Stop Recording** button is pressed.
- Once stopped, the voice input is saved as `temp_audio.wav`.
- This audio file is sent to `speech_to_text.py`, which processes the file using the **Vosk** offline speech recognition model. The `transcribe_audio` function extracts the textual content and returns it to `app.py`.
- 

## 2. Intent Recognition & Module Mapping

- `app.py` forwards the transcribed text to `main.py`.
- `main.py` utilizes the `intent_classifier.py` script, which applies **zero-shot classification** using the `facebook/bart-large-mnli` model to detect the intent behind the query.
- Based on the classified intent, `main.py` dynamically imports and invokes the corresponding module from the `modules/` directory.

Modules include:

- `get_net_income.py`, `get_revenue.py`, `get_cash_flow.py`, `get_balance_sheet.py`, etc.
- Each module connects to the **Financial Modeling Prep (FMP) API** to retrieve specific financial data and returns the result to `main.py`, where it's stored in an output dictionary.

### 3. RAG-Based Retrieval

- If the **RAG (Retrieval-Augmented Generation)** option is selected, the system further processes the query:
  - It uses **spaCy's en\_core\_web\_lg** model for **Named Entity Recognition (NER)** to identify key components like the organization name, financial metric, and timeline.
  - These are stored as key-value pairs and passed to the RAG pipeline.

#### RAG Mechanism:

- Uses **all-MiniLM-L6-v2** for embedding both query and dataset entries.
- Applies **cosine similarity** with **FAISS** to find relevant matches ( $\geq 75\%$  similarity).
- If a match is found, the system retrieves contextual data and passes it to the **Gemma 2B** LLM for generating a natural-language response.
- The final RAG-based response is sent back to **main.py** and stored alongside the FMP response.

### 4. Fallback Logic for Robustness

- If the FMP API fails, the RAG pipeline is automatically triggered to retrieve data from local CSV files, Excel sheets, or tabular datasets.
- If both the FMP API and RAG fail, the system invokes **web\_search.py**, which performs a **privacy-focused search using DuckDuckGo**, prioritizing results with numerical values for accuracy.
- Additionally, an **SQL query tool** (using **sql\_db.py**) allows structured querying from external data sources by converting them into an **SQLite** database.

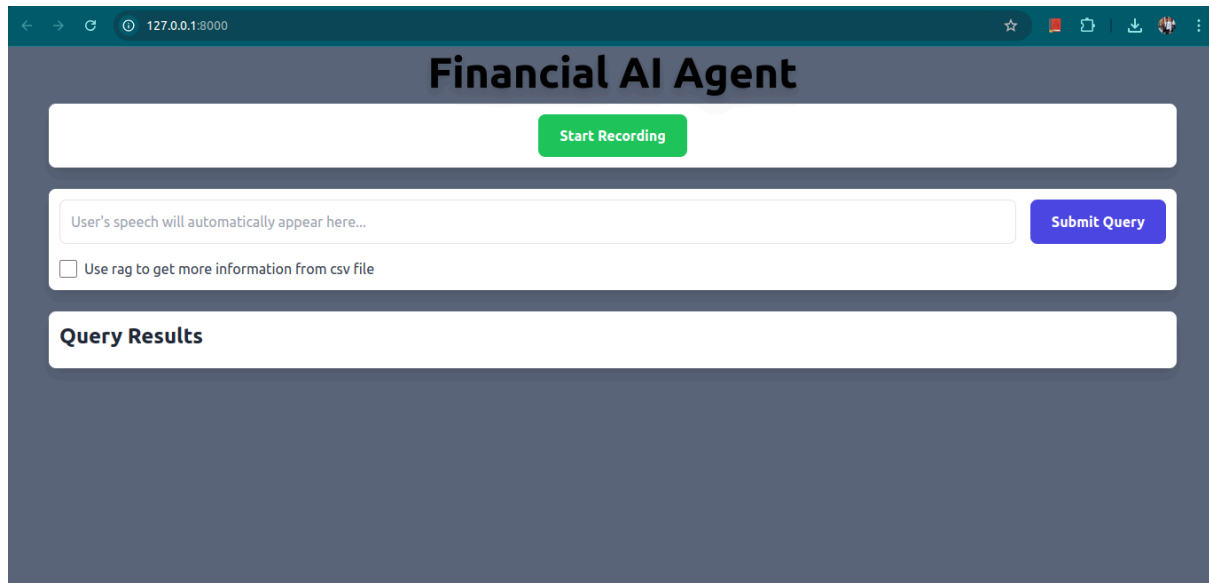
### 5. Error Handling

Robust error-handling mechanisms are built into every step of the process to ensure the system remains reliable and user-friendly during unexpected failures or incomplete inputs.

# User Interface Overview

**Start the Server:** Run `app.py` using:

```
```python
uvicorn app:app --reload
```
```



- Visit <http://127.0.0.1:8000/> in your browser.
- **Start Recording Button:**  
Begins recording the user's voice. Button toggles to **Stop Recording**.
- **Stop Recording Button:**  
Ends recording and saves the audio as a `.wav` file. Automatically sends it for speech-to-text processing.
- **Query Form:**  
Displays the transcribed text. Users can edit the query if needed.
- **Submit Query Button:**  
Sends the processed text to `main.py` for execution and retrieval.
- **RAG Checkbox:**  
If enabled, combines **RAG-based retrieval** with API responses for enhanced results.
- **Query Results Section:**  
Displays the final output from the API, RAG module, or web search.

# External Data Sources

## Financial Modeling Prep (FMP)

FMP is the open-source for the most reliable and accurate Stock Market API and Financial Data API. They provides real-time stock prices, financial statements, or historical data or a comprehensive solution to meet all the financial data needs for free.

**API key link:** <https://site.financialmodelingprep.com/developer/docs/dashboard>

N.B: They provide a free tier of 250 request per day.

### Sample Endpoints:

1. <https://financialmodelingprep.com/api/v3/income-statement/{ticker}>
2. <https://financialmodelingprep.com/api/v3/historical-price-full/{ticker}>
3. <https://financialmodelingprep.com/api/v3/balance-sheet-statement/{ticker}>

## Financial data sp500 companies.csv (kaggle)

This dataset represents 500 company's financial statements based on USA to analyze the company's stock price in the share market. The data is scrapped from the yahoo finance API.

### Dataset link:

<https://www.kaggle.com/datasets/pierrelouisdanieau/financial-data-sp500-companies/data>

### Data description:

Each line represents a financial report for a given date.

For each company there are 4 annual reports with 4 different dates:

- 2020-12-31
- 2021-03-31
- 2021-06-30
- 2021-09-30

The columns are :

1. firm : company name
2. Ticker : company ticker (the symbol)
3. Research Development
4. Income Before Tax
5. Net Income
6. Selling General
7. Administrative
8. Gross Profit

9. Ebit
10. Operating Income
11. Interest Expense
12. Income Tax Expense
13. Total Revenue
14. Total Operating Expenses
15. Cost Of Revenue
16. Total Other Income Expense Net
17. Net Income From Continuing Ops
18. Net Income Applicable To Common Shares

### Data preprocessing:

For retrieving the information process easier,

- Changed date time to year.

| Previous date                          | Changed into |
|--|--------------|
| 2021-03-31<br>2021-06-30<br>2021-09-30 | 2024         |
| 2020-12-31                             | 2023         |

- Changes column name.

| Previous Name                          | Current Name    |
|--|-----------------|
| Net Income                             | netIncome       |
| Gross Profit                           | netProfitMargin |
| Total Revenue                          | revenue         |
| Income Tax Expense                     | IncomeTax       |
| Interest Expense                       | InterestExpense |
| Research Development                   | Research        |
| Cost of Revenue                        | TotalCost       |
| Net Income Applicable to common shares | CommonShares    |



# Language Models

## facebook/bart-large-mnli

The `facebook/bart-large-mnli` model is a pre-trained transformer fine-tuned on the MNLI dataset, enabling it to perform **zero-shot text classification**. It excels at recognizing user intent by evaluating how well input text aligns with a set of candidate labels without needing task-specific training. This makes it ideal for applications like intent detection, topic tagging, and content categorization.

## en\_core\_web\_lg

`en_core_web_lg` is a large English NLP model provided by spaCy, featuring over 700k word vectors and pre-trained pipeline components for tasks like part-of-speech tagging, dependency parsing, and **named entity recognition**. It offers high accuracy and performance for analyzing complex English text. Its rich embeddings make it suitable for semantic similarity, intent recognition, and other advanced NLP tasks.

## all-MiniLM-L6-v2

`all-MiniLM-L6-v2` is a compact, efficient **sentence embedding model** from the SentenceTransformers library, based on a 6-layer MiniLM architecture. It is trained using a contrastive learning objective to produce high-quality embeddings for tasks like semantic search, clustering, and zero-shot classification. Despite its small size, it delivers strong performance and fast inference, making it ideal for real-time NLP applications.

## Ollama: Gemma 2B

The `gemma:2b` model is part of Google's Gemma family of lightweight, open-source language models developed by Google DeepMind. With approximately 2.5 billion parameters, it is designed for efficient performance across various NLP tasks. The model has been trained on a diverse dataset of web documents, including code and mathematical text, to handle a wide range of linguistic styles and topics .

## Dependencies

### ❖ Python Libraries:

- `vosk`, `pyaudio`, `pydub`: For speech recognition.
- `Spacy`, `spacy-lookups-data`: For name entity recognition.
- `pandas`: For CSV/Excel handling.
- `Langchain`, `faiss-cpu`, `sentence-transformers`, `fuzzywuzzy`, `sqlite3`
- `duckduckgo_search`: For web searches.
- `asyncio`, `httpx`: For asynchronous operations and HTTP requests.

- fastapi, uvicorn, jinja2: For building a local server api.
- ❖ System Packages:
  - Portaudio19-dev: For recording audio via web.
  - Ollama: For using cpu compatible opensource language models.

## Future Improvements

We can add new financial operation by creating a new Python script inside the `modules/` directory (e.g., `get_operating_expense.py`) and define a function that handles the API call or data processing logic. Then, update the **intent classifier's label list** (if needed) to recognize the new intent. The `main.py` script dynamically imports and executes modules based on the detected intent—so no changes to the core logic are required.

## Conclusion

The Financial AI Agent provides a robust solution for querying financial data using voice recognition, with a modular architecture that ensures flexibility and scalability. It can adapt any new functionalities related to financial operation. This project can fully run on a cpu machine ensuring lowest possible computation power is not a limitation.

# Testing and Integration

## API Testing

### 1. Get\_Financial\_Ratios:

The screenshot shows a web browser window with the address bar displaying '127.0.0.1:8000'. The page title is 'Financial AI Agent'. At the top, there is a green 'Start Recording' button. Below it is a text input field with the placeholder 'User's speech will automatically appear here...' and a blue 'Submit Query' button. A checkbox labeled 'Use rag to get more information from csv file' is present. The 'Query Results' section displays the following information:

- User\_Query:** what the investment ratio that microsoft allotted or research in last year
- Intent:** get\_financial\_ratios
- Entities:** {'ticker': 'MSFT', 'metric': 'Research', 'year': '2024', 'date': None}
- API Response:** MSFT's current ratio for 2024 is 1.27.
- Final Response:** MSFT's current ratio for 2024 is 1.27.

### 2. get\_stock\_price:

The screenshot shows the same 'Financial AI Agent' interface. The 'Query Results' section displays the following information:

- User\_Query:** the stock price of facebook
- Intent:** get\_stock\_price
- Entities:** {'ticker': 'META', 'metric': 'price', 'year': None, 'date': None}
- API Response:** META's current stock price is \$545.42.
- Retriever Response:** No relevant data found.
- Final Response:** META's current stock price is \$545.42. Additional Info found in the CSV: No relevant data found.

### 3. get\_market\_capitalization:

The screenshot shows a web browser window with the address bar displaying '127.0.0.1:8000'. The page title is 'Financial AI Agent'. At the top, there is a green 'Start Recording' button. Below it is a text input field with the placeholder 'User's speech will automatically appear here...' and a blue 'Submit Query' button. A checkbox labeled 'Use rag to get more information from csv file' is present. The 'Query Results' section displays the following information:

- User\_Query:** what is the current market value of apple
- Intent:** get\_market\_cap
- Entities:** {'ticker': 'AAPL', 'metric': 'mktCap', 'year': None, 'date': None}
- API Response:** AAPL's market cap is \$3079.83 billion.
- Retriever Response:** No relevant data found.
- Final Response:** AAPL's market cap is \$3079.83 billion. Additional Info found in the CSV: No relevant data found.

### 4. Get\_stock\_price:

The screenshot shows the same web browser window as above, but with the 'Query Results' section displaying information for a stock price query:

- User\_Query:** the stock price of facebook
- Intent:** get\_stock\_price
- Entities:** {'ticker': 'META', 'metric': 'price', 'year': None, 'date': None}
- API Response:** META's current stock price is \$545.42.
- Retriever Response:** No relevant data found.
- Final Response:** META's current stock price is \$545.42. Additional Info found in the CSV: No relevant data found.

## 5. Get\_cash\_flow:

The screenshot shows the 'Financial AI Agent' web interface. At the top, there's a browser address bar with '127.0.0.1:8000' and a notification 'You can paste the image from the clipboard.' The main header is 'Financial AI Agent'. Below it is a green 'Start Recording' button. A text input field contains 'User's speech will automatically appear here...' and a blue 'Submit Query' button is to its right. Below the input field is a checkbox labeled 'Use rag to get more information from csv file'. The 'Query Results' section displays the following information:

- User\_Query:** what is the cash flow of amazon last year
- Intent:** get\_cash\_flow
- Entities:** {'ticker': 'AMZN', 'metric': 'cashFlowFromOperatingActivities', 'year': '2024', 'date': None}
- API Response:** AMZN's cash from operations for 2024 is \$0.00 billion.
- Retriever Response:** No relevant data found.
- Final Response:** AMZN's cash from operations for 2024 is \$0.00 billion. Additional Info found in the CSV: No relevant data found.

## Rag Testing:

### 1. Get\_research\_info:

The screenshot shows the 'Financial AI Agent' web interface. At the top, there's a browser address bar with '127.0.0.1:8000' and a notification 'You can paste the image from the clipboard.' The main header is 'Financial AI Agent'. Below it is a green 'Start Recording' button. A text input field contains 'User's speech will automatically appear here...' and a blue 'Submit Query' button is to its right. Below the input field is a checkbox labeled 'Use rag to get more information from csv file'. The 'Query Results' section displays the following information:

- User\_Query:** the investment read that microsoft allotted for research in last year
- Intent:** get\_research\_info
- Entities:** {'ticker': 'MSFT', 'metric': 'Research', 'year': '2024', 'date': None}
- API Response:** Error fetching base response: GetResearchInfo.get\_data() got an unexpected keyword argument 'year' Using retriever to query CSV file...
- Retriever Response:** The investment read that Microsoft allotted for research in last year was \$5.6 billion.
- Final Response:** The investment read that Microsoft allotted for research in last year was \$5.6 billion.

At the bottom of the interface, there is a 'Terminal' button.

## 2. Get\_net\_income:

The screenshot shows the 'Financial AI Agent' interface. At the top, there's a 'Start Recording' button. Below it is a text input field with the placeholder 'User's speech will automatically appear here...' and a 'Submit Query' button. A checkbox labeled 'Use rag to get more information from csv file' is checked. The 'Query Results' section displays the following information:

- User\_Query:** what is the net income of tesla this year
- Intent:** get\_net\_income
- Entities:** {'ticker': 'TSLA', 'metric': 'netIncome', 'year': '2025', 'date': None}
- API Response:** TSLA's net income for 2025 is \$7.13 billion.
- Retriever Response:** The net income of Tesla in 2024 was \$43,800,000,000.
- Final Response:** TSLA's net income for 2025 is \$7.13 billion. Additional Info found in the CSV: The net income of Tesla in 2024 was \$43,800,000,000.

## 3. get\_income\_tax:

The screenshot shows the 'Financial AI Agent' interface. At the top, there's a 'Start Recording' button. Below it is a text input field with the placeholder 'User's speech will automatically appear here...' and a 'Submit Query' button. A checkbox labeled 'Use rag to get more information from csv file' is checked. The 'Query Results' section displays the following information:

- User\_Query:** income tax that walmart paid last year
- Intent:** get\_revenue
- Entities:** {'ticker': 'WMT', 'metric': 'netIncome', 'year': '2024', 'date': None}
- API Response:** WMT's revenue for 2024 is \$680.99 billion.
- Retriever Response:** The income tax paid by Walmart last year was approximately 4.27 billion dollars.
- Final Response:** WMT's revenue for 2024 is \$680.99 billion. Additional Info found in the CSV: The income tax paid by Walmart last year was approximately 4.27 billion dollars.

# Web Search Testing

## 1. get\_interest\_expense:

The screenshot shows a web application interface with a dark blue header. At the top, there is a green button labeled "Start Recording". Below this is a search bar with the placeholder text "User's speech will automatically appear here..." and a blue button labeled "Submit Query". A checkbox labeled "Use rag to get more information from csv file" is located below the search bar. The main content area is titled "Query Results" and contains the following information:

**User\_Query:** what are the total interest that meta paid to have bank  
**Intent:** get\_interest  
**Entities:** {'ticker': 'META', 'metric': 'InterestExpense', 'year': None, 'date': None}  
**API Response:** Error fetching base response: GetInterest.get\_data() got an unexpected keyword argument 'year' Using retriever to query CSV file...  
**Retriever Response:** No relevant data found.  
**Web Search Response:** Capitalized interest: Amount of interest costs capitalized disclosed as an adjusting item to interest costs incurred. Meta Platforms Inc. capitalized interest increased from 2022 to 2023 and from 2023 to 2024. Interest costs incurred: Total interest costs incurred during the period and either capitalized or charged against earnings.  
**Final Response:** Capitalized interest: Amount of interest costs capitalized disclosed as an adjusting item to interest costs incurred. Meta Platforms Inc. capitalized interest increased from 2022 to 2023 and from 2023 to 2024. Interest costs incurred: Total interest costs incurred during the period and either capitalized or charged against earnings.

## 2. get\_interest\_expense:

The screenshot shows a web application interface with a dark blue header. At the top, there is a green button labeled "Start Recording". Below this is a search bar with the placeholder text "User's speech will automatically appear here..." and a blue button labeled "Submit Query". A checkbox labeled "Use rag to get more information from csv file" is located below the search bar. The main content area is titled "Query Results" and contains the following information:

**User\_Query:** give a brief description of nvidia  
**Intent:** get\_interest  
**Entities:** {'ticker': 'NVDA', 'metric': None, 'year': None, 'date': None}  
**API Response:** Error fetching base response: GetInterest.get\_data() got an unexpected keyword argument 'year' Using retriever to query CSV file...  
**Retriever Response:** No relevant data found.  
**Web Search Response:** NVIDIA in Brief NVIDIA pioneered accelerated computing to tackle challenges no one else can solve. Our work in AI and digital twins is transforming the world's largest industries and profoundly impacting society. Learn more. Company History Founded in 1993, NVIDIA is the world leader in accelerated computing. Our  
**Final Response:** NVIDIA in Brief NVIDIA pioneered accelerated computing to tackle challenges no one else can solve. Our work in AI and digital twins is transforming the world's largest industries and profoundly impacting society. Learn more. Company History Founded in 1993, NVIDIA is the world leader in accelerated computing. Our