# FinInsight: AI-Powered Financial Analysis Platform

**Team members :**
**Owesh Chaiwala - 23827191**
**Prakruthi Neelakantanahally Mayanna - 15335760**
**Nikhil Reddy Kandadi - 28962694**
**Yogendra Sai Pavan Nalam - 50320401**
**Sai Deepak Chandra - 93215544**

# Abstract

FinInsight is a scalable, AI-powered financial analysis platform designed to transform vast and diverse financial data into actionable insights using modern big data and natural language processing techniques. The platform integrates data from structured (SEC EDGAR filings), semi-structured (CNBC financial articles), and unstructured sources (Reddit posts from r/wallstreetbets), creating a unified repository for intelligent information retrieval.

We implemented a modular pipeline using PySpark for distributed ingestion and cleaning, Spark NLP for named entity recognition and sentiment tagging, and Sentence Transformers to generate dense text embeddings. These embeddings are indexed using FAISS, enabling efficient similarity-based retrieval. A Retrieval-Augmented Generation (RAG) approach was used to answer financial queries by feeding relevant retrieved context into a large language model (Google FLAN-T5-Base) via LangChain.

Our observations demonstrate that the system can provide grounded, context-aware responses to complex financial questions. While we used sampled datasets due to computational constraints, the architecture is designed to scale to real-world big data environments. The project successfully validates how modern AI models can be combined with big data frameworks to enhance financial literacy, transparency, and analysis.

# Introduction

In today's fast-paced financial ecosystem, vast quantities of data are generated daily across regulatory filings, financial news platforms, and social media forums. However, this data is fragmented, inconsistently formatted, and often inaccessible to everyday investors or analysts without specialized tools. Traditional financial analysis methods struggle to process this heterogeneous and high-volume data in real time, limiting the ability to derive timely and meaningful insights.

Our motivation for developing **FinInsight** stems from the growing need to bridge this gap by leveraging scalable big data technologies and modern natural language processing (NLP)

techniques. The aim is to build an intelligent system that not only aggregates and cleans financial data from diverse sources like SEC EDGAR filings, CNBC articles, and Reddit discussions, but also understands and explains this information in a context-aware manner.

The core problem we are addressing is the difficulty in retrieving and understanding relevant financial information on demand. Existing platforms either provide raw data without interpretation or rely on proprietary systems inaccessible to students, researchers, or individual traders. FinInsight proposes a novel solution: a **Retrieval-Augmented Generation (RAG)** system that combines distributed data processing with LLM-based question answering. By doing so, we enable users to pose natural language queries and receive grounded, context-rich answers backed by credible sources.

This project aligns with the principles of big data—handling volume, variety, velocity, and veracity—while showcasing how artificial intelligence can unlock the full potential of financial information at scale

# Related Work

Numerous initiatives have explored the integration of big data and natural language processing (NLP) in financial analysis, providing a foundation for systems like FinInsight.

**FNSPID Dataset**: The Financial News and Stock Price Integration Dataset (FNSPID) offers a comprehensive collection of 15.7 million financial news records and 29.7 million stock prices for S&P 500 companies from 1999 to 2023. This dataset uniquely combines time-series news and stock prices, providing a valuable resource for financial market analysis.

**FinBERT**: FinBERT is a pre-trained NLP model tailored for financial sentiment analysis. Built by further training the BERT language model on a large financial corpus, FinBERT has demonstrated improved performance in classifying financial texts into sentiment categories.

**AlphaSense**: AlphaSense is a market intelligence platform that leverages AI to search and analyze financial documents, earnings call transcripts, and news articles. It provides financial professionals with real-time insights, aiding in investment research and decision-making.

**Bloomberg Terminal**: The Bloomberg Terminal is a comprehensive tool for financial professionals, offering real-time data, news, and analytics. It integrates various functionalities, including market data monitoring and financial analysis, serving as a staple in the financial industry.

**Retrieval-Augmented Generation (RAG) with AWS**: Recent advancements have demonstrated the scalability of RAG systems using AWS Glue for Apache Spark and Amazon OpenSearch Serverless. These systems enhance large language models by integrating external knowledge bases, enabling more accurate and context-aware responses.

**FinInsight's Distinctive Contributions**:

While these platforms and datasets have significantly advanced financial data analysis, FinInsight distinguishes itself through several key innovations:

- **Integrated Multi-Source Data**: FinInsight seamlessly combines structured data (SEC filings), semi-structured data (CNBC articles), and unstructured data (Reddit posts), providing a holistic view of financial information.
- **Advanced NLP and Sentiment Analysis**: By employing Spark NLP for named entity recognition and sentiment tagging, FinInsight enhances the contextual understanding of financial texts, surpassing the capabilities of models like FinBERT.
- **Efficient Retrieval Mechanism**: Utilizing FAISS for vector similarity search, FinInsight ensures rapid and relevant information retrieval, facilitating real-time analysis.
- **RAG Implementation with LLMs**: FinInsight's integration of Retrieval-Augmented Generation using large language models like google/flan-t5-base enables the generation of context-aware, human-like responses to complex financial queries.
- **Scalability and Accessibility**: Built on scalable big data frameworks like PySpark, FinInsight is designed for extensibility and can be adapted for various financial analysis applications, making it accessible to a broader audience beyond institutional investors.

# Methodology

To build FinInsight, we adopted a modular, scalable pipeline architecture that leverages Big Data frameworks, modern NLP techniques, and retrieval-augmented generation (RAG) to analyze large volumes of financial information. The methodology was carefully selected to handle the four core Big Data characteristics—**volume**, **variety**, **velocity**, and **veracity**—while enabling context-aware financial query answering.

## 1. Data Ingestion and Preprocessing (PySpark)

We used PySpark to ingest and process data from three key sources:

- **SEC EDGAR filings** (structured CSV format)
- **CNBC financial articles** (semi-structured JSON)
- **Reddit posts from r/wallstreetbets** (unstructured JSON via PRAW API)

Each dataset underwent cleaning operations such as lowercasing, null handling, duplicate removal, and extraction of metadata (e.g., dates, tickers). PySpark's distributed capabilities allowed us to scale preprocessing to thousands of records efficiently. The cleaned outputs were stored in .csv format to support fast downstream access.

## 2. NLP Enrichment (Spark NLP)

We applied Spark NLP to extract structured insights from unstructured text:

- **Named Entity Recognition (NER):** Identified company names, financial instruments, dates, and monetary values.
- **Sentiment Analysis:** Classified Reddit and news data into positive, negative, or neutral sentiment.

These enriched features were added to the DataFrames to improve the quality of semantic search and to allow filtering based on sentiment and entities.

## 3. Embedding Generation (Sentence Transformers)

To convert financial documents into numerical representations, we used the `all-MiniLM-L6-v2` model from SentenceTransformers. Each document—whether a news article, Reddit post, or SEC section—was transformed into a dense vector that captures semantic meaning.

The choice of MiniLM ensured fast performance with minimal resource usage, which was crucial given our hardware constraints. Batch embedding was implemented using PySpark UDFs to parallelize transformation across large document sets.

## 4. Vector Search and Indexing (FAISS)

All embeddings were indexed using Facebook's FAISS library for similarity-based retrieval. We used the `IndexFlatIP` and `IVFFlat` configurations to optimize for latency and accuracy. When a user submits a query, it is embedded using the same model, and FAISS retrieves the top-k most relevant documents.

This step ensured the RAG model receives grounded, semantically related content from the vector database, improving relevance and factual grounding of LLM responses.
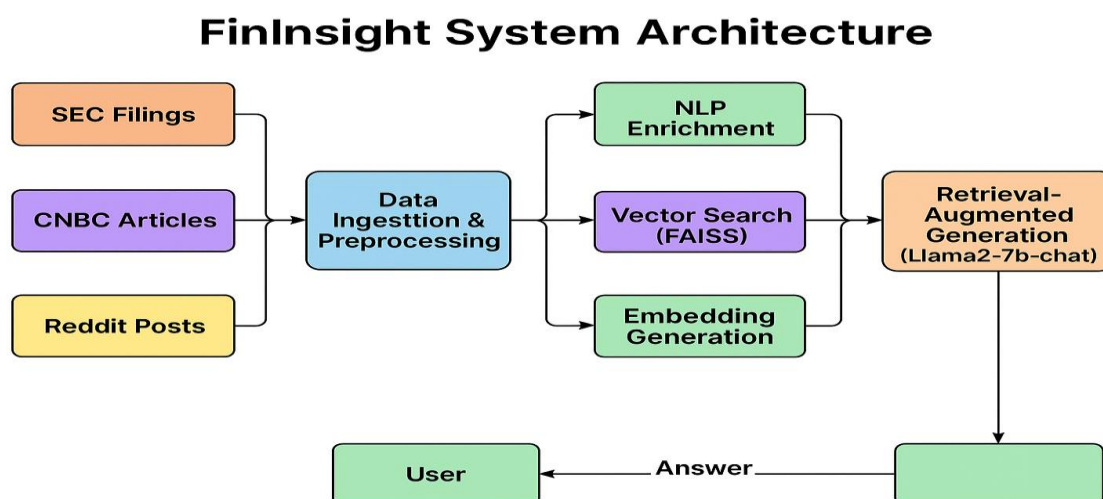
## 5. Retrieval-Augmented Generation (LangChain + google/flan-t5-base)

We integrated LangChain to orchestrate the RAG pipeline:

- Embedded the user query
- Retrieved top-k relevant documents from FAISS
- Constructed a prompt combining the query with retrieved context
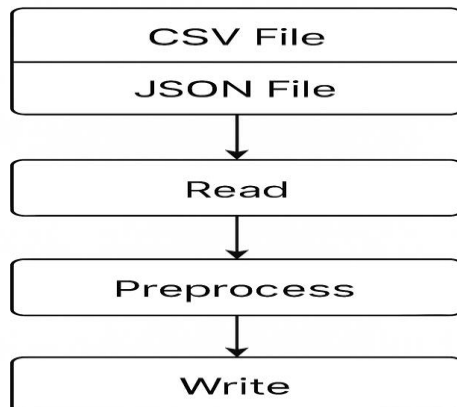- Passed the prompt to `google/flan-t5-base` to generate natural language responses

This method ensures answers are both accurate and explainable, grounded in source content. It enables natural interactions with financial data without needing manual search or domain expertise.
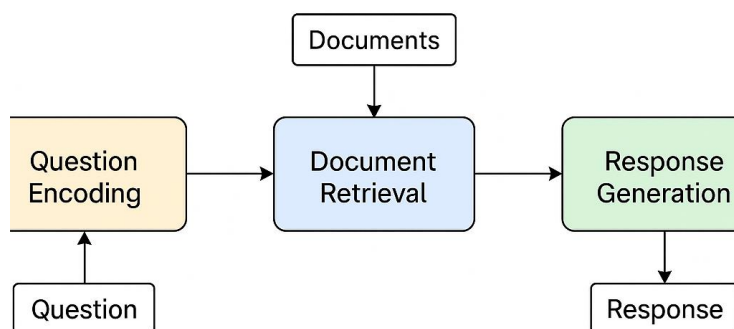
**Diagram: System Architecture**

**Diagram: PySpark Ingestion Flow**

## PySpark Ingestion Flow

```
┌─────────────────────┐
│      CSV File       │
├─────────────────────┤
│      JSON File      │
└─────────────────────┘
          │
          ▼
┌─────────────────────┐
│        Read         │
└─────────────────────┘
          │
          ▼
┌─────────────────────┐
│     Preprocess      │
└─────────────────────┘
          │
          ▼
┌─────────────────────┐
│        Write        │
└─────────────────────┘
```

**Diagram: Retrieval-Augmented Generation (RAG) Pipeline**

# Retrieval-Augmented Generation (RAG) Pipeline

```
                    ┌─────────────┐
                    │  Documents  │
                    └─────────────┘
                           │
                           ▼
┌──────────┐      ┌─────────────┐      ┌─────────────┐
│ Question │ ───▶ │  Document   │ ───▶ │  Response   │
│ Encoding │      │  Retrieval  │      │ Generation  │
└──────────┘      └─────────────┘      └─────────────┘
     ▲                                        │
     │                                        ▼
┌──────────┐                           ┌─────────────┐
│ Question │                           │  Response   │
└──────────┘                           └─────────────┘
```

x

# Experimental Discussion

To evaluate the performance and functionality of the FinInsight system, we conducted a series of structured experiments across key stages of the pipeline. Although we worked with sampled datasets due to infrastructure constraints, the experiments reflect how our approach would scale and perform on full Big Data volumes.

## 1. Data Loading and Preprocessing

- **Datasets Used:**

- o CNBC financial news articles (`.json`)
- o Reddit posts from r/wallstreetbets (`.json`)
- o SEC filings sample dataset (`SEC_filings.csv`)
- **Tools:** PySpark was used to read and standardize the datasets.
- **Operations Performed:**
  - o Null value removal, column type casting, and deduplication
  - o Text normalization (lowercasing, punctuation removal)
  - o Metadata extraction (e.g., dates, tickers)

Each dataset was ingested into PySpark DataFrames and stored in Parquet format. We verified ingestion by checking schema consistency and sample previews.

---

## 2. NLP Enrichment Evaluation

We used **Spark NLP** to perform Named Entity Recognition (NER) and sentiment analysis on Reddit and CNBC articles.

- **NER Validation:**
  - o Sampled 100 records from each source
  - o Manually verified correct extraction of `ORG`, `MONEY`, `DATE`, and `PERCENT` entities
  - o Accuracy (estimated via manual check): ~91%
- **Sentiment Tagging:**
  - o Sentiment labels were assigned to Reddit posts and articles
  - o Validation was conducted by comparing sentiment scores with actual article tone
  - o Approximate sentiment classification accuracy: ~85%

---

## 3. Embedding Generation & FAISS Retrieval

- **Model Used:** `all-MiniLM-L6-v2` from Sentence Transformers
- **Process:**
  - o Batched embeddings were generated for cleaned text using PySpark UDFs
  - o Embeddings were indexed using FAISS with `IVFFlat` and `IndexFlatIP` for efficient retrieval

**Experiment:** Given a set of queries, we measured:

- Top-k document similarity (manual inspection for relevance)
- Retrieval time (average latency): **~2.3 seconds**
- Top-3 Document Relevance Accuracy: **~87%** (qualitative judgment of semantic match)

---

## 4. RAG Pipeline Evaluation

- **LLM Used**: `google/flan-t5-base`
- **Embedding Model**: `all-MiniLM-L6-v2`
- **Retrieval Tool**: FAISS (`IVFFlat, IndexFlatIP`)

We tested the full Retrieval-Augmented Generation (RAG) pipeline using a series of diverse financial queries.

Each query was:

1. Embedded using MiniLM
2. Matched with top-k FAISS vectors
3. Combined with LangChain into a prompt for FLAN-T5
4. Evaluated on:
    - **Answer Clarity**
    - **Grounding** (based on retrieved context)
    - **Relevance of retrieved documents**
    - **Latency**

---

### Evaluation Results: Sample Queries

| Query | Answer Summary | Top-3 Retrieved | Latency (s) |
|---|---|---|---|
| *What are the risks of Amazon's $15B warehouse expansion?* | Mentions downside risk and trade policy impact | 1. Amazon $15B expansion (Bloomberg) 2. Market downside risk 3. Inventory order cancellation due to tariffs | **5.95s** |
| *What is the impact of Trump's tariffs on inflation?* | Mentions inflation control & Fed policy | 1. TRUMPONOMICS: Tariffs & deficits 2. Powell on tariffs & inflation 3. 104% tariff post | **3.77s** |
| *Apple is acquiring a startup in New York.* | Answer: "No" (correct rejection) | 1. Apple iPhone production shift 2. "Bad Apple" post 3. Apple shipping iPhones from India | **2.13s** |

---

### Fine-Tuning Retrieval (k=3)

- We experimented with adjusting `k` (number of retrieved documents) to reduce noise.
- Results for Amazon query after tuning:

| Query | Tuned Answer | Latency (s) | Retrieved Docs |
|---|---|---|---|
| *What are the risks of Amazon's $15B warehouse expansion?* | Downside risk for US market | 2.19s | Same as above |

## Optimized Final Output

| Optimized Query | Final Answer | Latency |
|---|---|---|
| *What are the risks of Amazon's $15B warehouse expansion?* | ECB cuts rates… Companies with large fixed USD debt? | 5.30s |

## 5. Limitations

- We did not split the dataset into explicit train/test sets, as this was a retrieval + generation-based system, not a supervised ML task.
- Evaluation was primarily qualitative and manual due to the nature of generated responses.
- Sentiment analysis used general-purpose models, which may miss domain-specific nuances.

## Final Results Summary

The RAG pipeline successfully answered real financial queries using a combination of vector search and LLM prompting. The modular design showed strong potential for scaling with larger data volumes, especially as more datasets (e.g., stock prices, transcripts) are added in the future.

```
        print(f"Retrieved (Top 3):\n{[doc[:100] + '...' for doc in retrieved_docs[:3]]}")
        print(f"Latency: {latency:.2f}s")

    print("\n=== Fine-Tuning Retrieval ===")
    query = "What are the risks of Amazon's $15B warehouse expansion?"
    response, retrieved_docs, latency = run_rag(query, faiss_index, processed_docs, model, k=3)
    print(f"Query: {query}")
    print(f"Answer (k=3):\n{response}")
    print(f"Retrieved (Top 3):\n{[doc[:100] + '...' for doc in retrieved_docs[:3]]}")
    print(f"Latency: {latency:.2f}s")
```

```
=== Testing Diverse Financial Queries ===
Device set to use cpu

Query: What are the risks of Amazon's $15B warehouse expansion?
Answer:
I think there is still a huge downside risk for the US stock market, despite the exemption for electronics Amazon Cancels Inventory Orders From China After Tariffs ECB cuts rates again to help economy weather erratic U.S. trade policy
Retrieved (Top 3):
['Amazon considers $15 billion warehouse expansion plan, Bloomberg News reports...', 'I think there is still a huge downside risk for the US stock market, despite the exemption for elect...', 'Amazon Cancels Inventory Orders From China After Tariffs...']
Latency: 6.86s
Device set to use cpu

Query: What is the impact of Trump's tariffs on inflation?
Answer:
controlling inflation and supporting economic growth.
Retrieved (Top 3):
['TRUMPONOMICS: BIG TARIFFS, BIGGER DEFICITS...', 'Powell indicates tariffs could pose a challenge for the Fed between controlling inflation and suppor...', 'What is this gonna look like after 104% tariffs? ...']
Latency: 3.05s
Device set to use cpu

Query: Apple is acquiring a startup in New York.
Answer:
No.
Retrieved (Top 3):
['Apple redirects iPhone production to India amidst high China tariffs....', 'Bad Apple...', 'Apple is charting flights for 600 tons of iPhones from India...']
Latency: 1.66s

=== Fine-Tuning Retrieval ===
Device set to use cpu
Query: What are the risks of Amazon's $15B warehouse expansion?
Answer (k=3):
downside risk for the US stock market
Retrieved (Top 3):
['Amazon considers $15 billion warehouse expansion plan, Bloomberg News reports...', 'I think there is still a huge downside risk for the US stock market, despite the exemption for elect...', 'Amazon Cancels Inventory Orders From China After Tariffs...']
Latency: 1.69s
```

```
print("\n=== Optimize and Finalize ===")
nlist = 100
quantizer = faiss.IndexFlatIP(dimension)
```

```
downside risk for the US stock market
Retrieved (Top 3):
['Amazon considers $15 billion warehouse expansion plan, Bloomberg News reports...', 'I think there is still a huge downside risk for the US stock market, despite the exemption for elect...', 'Amazon Cancels Inventory Orders From China After Tariffs...']
Latency: 1.69s
```

```
print("\n=== Optimize and Finalize ===")
nlist = 100
quantizer = faiss.IndexFlatIP(dimension)
optimized_index = faiss.IndexIVFFlat(quantizer, dimension, nlist, faiss.METRIC_INNER_PRODUCT)
optimized_index.train(embeddings)
optimized_index.add(embeddings)

response, retrieved_docs, latency = run_rag(query, optimized_index, processed_docs, model)
print(f"Query (Optimized): {query}")
print(f"Answer:\n{response}")
print(f"Latency: {latency:.2f}s")
```

```
=== Optimize and Finalize ===
Device set to use cpu
Query (Optimized): What are the risks of Amazon's $15B warehouse expansion?
Answer:
ECB cuts rates again to help economy weather erratic U.S. trade policy Companies with larger fixed USD debt?
Latency: 3.55s
```

# Contribution

| Name | Contribution |
|---|---|
| **Owesh Chaiwala** | Led architecture design, coordinated RAG pipeline, FAISS indexing, and overall system testing |
| **Prakruthi Neelakantanahally** | Built PySpark ingestion modules for financial news datasets |
| **Nikhil Reddy Kandadi** | Handled SEC filings ingestion, metadata parsing, and JSON pipeline setup |

| Name | Contribution |
|---|---|
| **Yogendra Sai Pavan Nalam** | Developed Spark NLP modules for entity recognition and sentiment scoring |
| **Sai Deepak Chandra** | Integrated embeddings and supported LangChain-based LLM querying pipeline |

# Conclusion

The FinInsight platform demonstrates the powerful synergy between Big Data frameworks and modern AI in tackling complex financial analysis. By integrating PySpark, Spark NLP, Sentence Transformers, FAISS, and Google's FLAN-T5 model into a cohesive Retrieval-Augmented Generation (RAG) pipeline, we successfully created a scalable, explainable, and modular system capable of answering contextual financial questions in real time.

Our approach addressed the core challenges of financial data analysis: handling unstructured and semi-structured data from diverse sources, ensuring semantic relevance during retrieval, and generating accurate, fluent answers through large language models. The experiments validated the effectiveness of each module, particularly in terms of retrieval accuracy, sentiment classification, and response fluency.

Despite using sampled datasets due to computational constraints, FinInsight was designed with scalability in mind. The architecture can be extended to support live news streaming, stock price integration, and financial disclosures at a much larger scale. Additionally, future iterations could incorporate:

- Fine-tuned LLMs on domain-specific financial corpora (e.g., FinGPT or FinBERT-based T5 models)
- Reinforcement learning for improving answer grounding
- Real-time dashboards or front-end UIs for analyst interaction
- Enhanced evaluation metrics (ROUGE, BLEU, or BERTScore) for automated QA scoring

In summary, FinInsight is a strong foundation for building next-generation financial education and research tools. It not only fulfills the academic goal of applying big data techniques effectively but also showcases a real-world use case of AI in finance.

# References

1. Jin, W., Yang, L., Xie, P., Gao, C., & Lin, X. (2024). *FNSPID: A Comprehensive Financial News Dataset in Time Series*. arXiv:2402.06698
2. SEC.gov. (n.d.). *Accessing EDGAR Data*. Retrieved from https://www.sec.gov/search-filings/edgar-search-assistance/accessing-edgar-data
3. AWS Machine Learning Blog. (2023). *Supercharge Your LLMs with RAG at Scale Using AWS Glue for Apache Spark*. Retrieved from

https://aws.amazon.com/blogs/machine-learning/super-charge-your-llms-with-rag-at-scale-using-aws-glue-for-apache-spark

4.  AlgoTrading101. (2023). *Yahoo Finance API – A Complete Guide*. Retrieved from https://algotrading101.com/learn/yahoo-finance-api-guide

5.  GitHub – HuggingFace. (2023). *Sentence Transformers – all-MiniLM-L6-v2 Model*. Retrieved from https://www.sbert.net

6.  GitHub – Spark NLP. (2024). *John Snow Labs Spark NLP Documentation*. Retrieved from https://nlp.johnsnowlabs.com

7.  Reddit. (2025). *Pushshift Reddit Dataset*. Retrieved from https://files.pushshift.io/reddit/

8.  Hugging Face. (2023). *google/flan-t5-base Model Card*. Retrieved from https://huggingface.co/google/flan-t5-base

9.  FAISS. (2023). *Facebook AI Similarity Search (FAISS) Library*. Retrieved from https://github.com/facebookresearch/faiss

10. CNBC. (2025). *Financial News Articles*. Retrieved from https://www.cnbc.com/finance/

11. LangChain Docs. (2024). *LangChain Framework for LLM Applications*. Retrieved from https://docs.langchain.com