# An End-to-End Data Analytics Pipeline - Stock Performance Visualization

*By Anshika Goel, Sai Chaitanya Munagala, San Jose State University*
*Applied Data Science Department*
*22/04/2025*

## Abstract

**Building scalable and efficient data pipelines is essential for analytics and business intelligence in the era of data-driven industries. This project explicitly aims and demonstrates an end-to-end ELT [Extract, Load, Transform] pipeline architecture using various tech stacks like dbt, Apache Airflow, Looker, and Snowflake. The ETL pipeline extracts historical stock price data from Apple, Google, and Microsoft using yFinance, loads it into Snowflake, transforms the data using dbt models, and presents insights through Looker visualizations. This report provides an in-depth technical breakdown of each component and how they interact to form a cohesive and modular analytical system. It showcases the importance of automation in the data world and generates quick, easy efficient insights to the business needs that transform the business decisions.**

## 1. Introduction

### A. Background

The vast exponential growth from various sources has created an urge for robust and efficient scalable engineering workflows. Traditionally, ETL systems face challenges when dealing with the enormous amount of data, popularly called big data, and complex schemas. The modern ELT approach is the solution, it leverages cloud data warehouses for scalability and performance. In this project, we aim to design and deploy an automated ELT approach that not only ingests data efficiently but can also perform some transformation close to the storage layer. This minimizes I/O bottlenecks and supports efficient downstream analytics. This shift also aligns with the growing demand for near-real-time insights and reduced data latency in decision-making. By leveraging tools like dbt and Airflow, ELT workflows become more modular, testable, and maintainable. Such a setup always lays good groundwork and foundation for advanced analytical use cases.

### B. Problem Statement

In the current data-driven world, almost every organization deals with disparate datasets that need to be ingested, cleaned, enriched, and visualized in near real-time. The absence of automated orchestration, scalable layers, manual transformations, and integrated reporting tools not only leads to manual errors but also results in slow reporting cycles. The lack of visibility is also an issue that needs to be addressed. This project aims to challenge the automation of stock market data analytics using open-source orchestration tools and scalable cloud data platforms. In the end, it also visualises the performance obtained from the transformation to see the insights.

### C. Objectives

1. Extract stock data for Apple, Google, and Microsoft for the last 180 days using the *Yahoo Finance API*.
2. Build a modular and maintainable *ELT* pipeline using *Apache Airflow*.
3. Use *dbt* to define, test, and version control transformation logic.
4. Visualize meaningful insights via *Looker*.
5. Timely demonstrates *reproducibility, modularity,* and *extensibility* in the pipeline.

**D. Requirements and Specifications**

This project requires leveraging Apache Airflow, dbt, Snowflake, and Looker. The system is required to extract, load, and transform stock price data and visualize key metrics using a BI tool. The technical requirements and specifications for these are classified as below:

1. Airflow: Two DAGS, one for ETL and the other for ELT, to trigger dbt transformations.
2. dbt: Models created for metric calculations. Profiles.yml file for Snowflake connectivity.
3. Snowflake: A Warehouse and 2 Schemas for populating tables.
4. Looker: For visualizing data and drawing patterns as per the business requirements.

---

## 2. System Architecture

**A. Technology Stack**

1. *Apache Airflow*: A Workflow orchestration tool to manage and automate data pipelines.
2. *Snowflake*: Cloud-based data warehouse for storing and processing raw and curated data.
3. *dbt* [Data Build Tool]: Transformation tool that manages completely through *SQL*-based transformations with lineage, testing, and documentation.
4. *Looker*: A Data visualization platform to create dashboards and charts from Snowflake tables.
5. *yFinance*: Open-source Python library to fetch historical stock prices and indicators.
6. *Docker*: Containerization technology used to encapsulate the entire stack.
7. *Python:* Scripting language used for *DAG* and API integration.



**B. Architecture Overview**

The architecture follows an ELT pattern:

1. *Extract:* yFinance fetches historical data.
2. *Load:* Data is loaded into Snowflake's schema using Snowflake Connector.
3. *Transform:* dbt runs transformation models on top of tables to create analytics-ready tables in the analytics schema.
4. Visualize: *Looker* connects to Snowflake and visualizes curated tables.



This architecture keeps the pipeline modular, scalable and easy to maintain. Each of the data extraction, transformation and visualization can be monitored and improved independently. Tools like dbt and Looker make it easier to collaborate and deliver insights faster. This setup ensures the analytics layer always reflects the latest, most reliable data
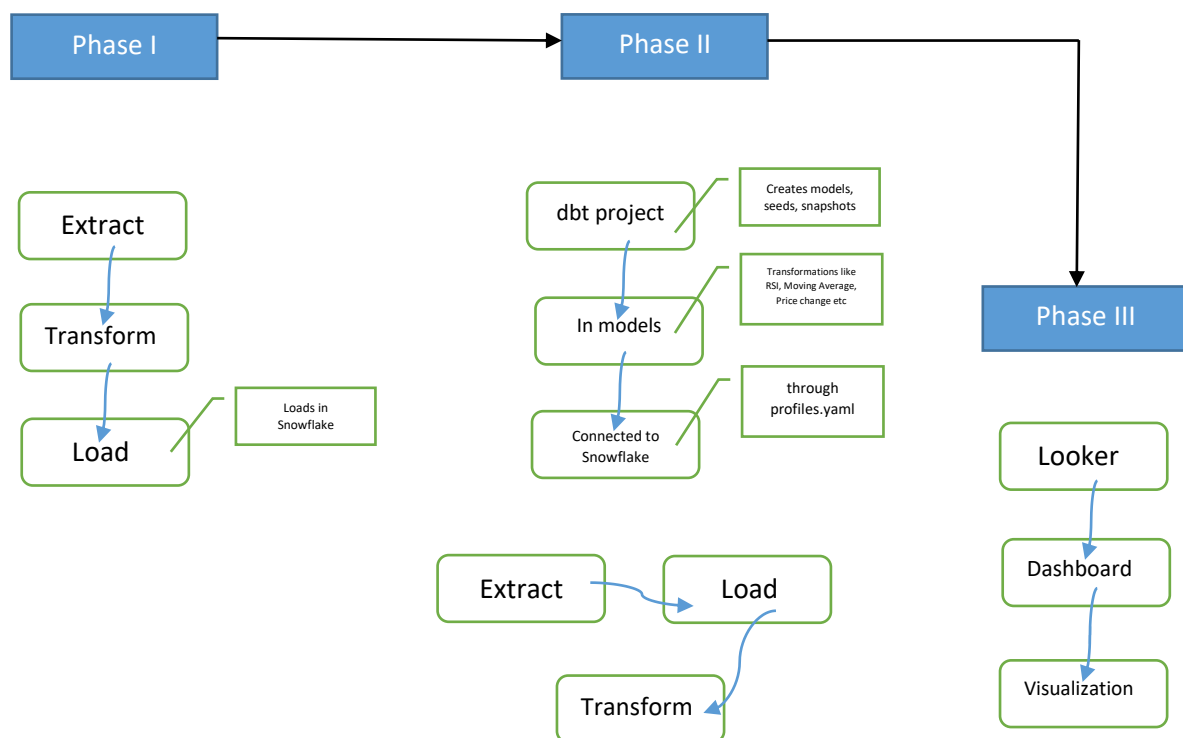
## 3. Implementation Phase and Overall System diagram

**A. Implementation Phase**

The implementation of the project was executed in three core stages: They are:

1. *ETL (Extract, Transform, Load)* using *Apache Airflow*,
2. *ELT modelling* using *dbt (data build tool)*, and the automation through *ELT* DAG.
3. *Data visualization* using *Looker*.

Each of these above stages is strategically planned and well implemented to ensure smooth, efficient data flow and to maintain transformation accuracy, and also for analytical insight generation. Below is a concise narrative of how the implementation phase unfolded at each of its respective stages.

**B. Overall System Diagram:**



First, to begin with, the environment has to be set up; the entire infrastructure for running the workflow was initially set up using Docker. A pre-configured Airflow Docker setup was utilized, where services like the Airflow webserver, scheduler, worker, and triggerer were defined within the *docker-compose.yml* file. The file was also updated to ensure all the necessary Python packages, such as snowflake-connector-python, dbt-snowflake, and apache-airflow-providers-snowflake, and at a later stage, Looker, a powerful dashboard for Business Intelligence, is required for installation. This Looker setup will later help us connect with Snowflake at a later stage to run visuals and draw insights. Once this was in place, the Airflow stack was launched using the *docker compose up* command.

This setup provides a solid foundation to orchestrate our pipeline in a clean and efficient way. The above system diagram at every necessary layer ensures SQL transactions with a try/except method, validates tests at every dbt transformation stage, and checks by maintaining idempotency at every required step.
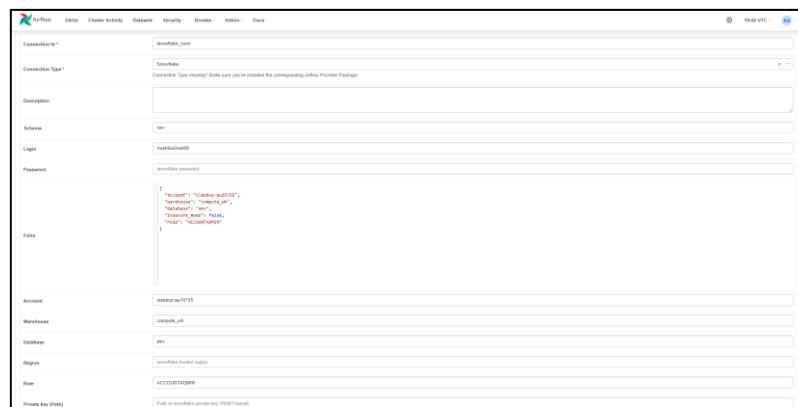
## 4. Methodology

As discussed in the Implementation phase, the methodology will be categorised into three crucial phases, namely, Phase I, Phase II, and Phase III

### A. Phase I

The first phase of data processing focused on ETL. A custom DAG *(StockPrice_DataLoad_Pipeline.py)* was created to automate the extraction of stock price data by using the *yfinance*. The DAG was configured to fetch the last 180 days of data for selected stock symbols [i.e., Apple, Google, Microsoft]. This data was fetched using a PythonOperator through three tasks, namely, extract, transform, and load within the DAG, and transformed into a structured format suitable for Snowflake.



[i] Screenshot Airflow ETL - DAG with its respective tasks



[i] Screenshot Airflow Connections

Subsequently, this data was inserted into Snowflake. Before insertion, a Snowflake connection was established within the *Airflow UI*, supplying account credentials, *database, schema, warehouse,* and *role*. This connection enabled the DAG to push data securely and seamlessly from Airflow into the Snowflake data warehouse. This is now officially called a *data pipeline* focused on *ETL*, insured best with SQL transactions and idempotency checks at every necessary stage.

### B. Phase II

The Second phase, an important phase where the project will go through its transition, called the ELT phase, using dbt. For this, a fresh DBT project was initialized locally, which generates the necessary folder structure and configuration files. The connection to Snowflake was defined in *profiles.yml*, where authentication parameters like database names and other metadata were provided. With the dbt project configured, model files were written under the */models directory*. These SQL-based models transform the raw stock data into analytical insights. For instance, *stock_summary.sql* was designed to calculate technical indicators like the moving averages,

the Relative Strength Index, the Simple moving averages, gain/loss ratios, and more. These transformations help to derive cleaner and more insightful data for analysis.

To execute the dbt transformations within the pipeline, a second Airflow DAG *(BuiltELT_dbt.py)* was created. This DAG used a BashOperator to run dbt run directly from the Airflow scheduler container. An alternate manual option to trigger these transformations is to prompt *dbt run*, *dbt test* through command prompts, but since this project focuses completely on automation without manual intervention, we choose to create a Dag *(BuiltELT_dbt.py)* as mentioned earlier, which primarily focuses on triggering the transformations. When triggered, this DAG compiles and runs all *SQL* models defined in the dbt project, thereby materializing new analytical tables into the *analytics* schema of Snowflake. This represents the complete automation of the ELT process from transformation design to execution. This, now from here, will be called a complete *end-to-end ELT Automated Pipeline.*



[i] Screenshot Airflow ELT DAG and its respective tasks

**[i] Detailed Table Structure:**

This is the output table generated via dbt transformations. It is built on top of raw stock data and includes analytical metrics such as RSI and Moving Averages, etc.

| Field Name | Data Type | Description | Constraints |
|---|---|---|---|
| **Date** | Date | Date of Stock Metrics | Not Null, Primary Key along with Symbol |
| **Symbol** | String | Ticker symbol of the stock | Not Null, Primary Key, along with Date |
| **Open** | Float | Price of the stock at the beginning | Nullable |
| **Close** | Float | Price of the stock on the closing day | Nullable |
| **High** | Float | Highest Price of the day | Nullable |
| **Low** | Float | Lowest Price of the day | Nullable |
| **Volume** | Float | Size of the stock | Nullable |
| **Previous close** | Float | Previous day closing Price | Nullable |
| **Change** | Float | [Closing Price – Previous price] | Nullable |
| **Gain** | Float | Positive change in price, else 0 | Nullable |
| **Loss** | Float | Negative change in price, else 0 | Nullable |
| **Avg_Gain** | Float | 15-day rolling average of gains | Nullable |
| **Avg_Loss** | Float | 15-day rolling average of losses | Nullable |
| **RSI** | Float | 15-day RSI with avg_gain and avg_loss | Nullable |
| **Moving_avg** | Float | 20-day rolling average Closing price | Nullable |

## E. Phase III

After building and loading all analytical models into *Snowflake*, the final phase involves connecting the data warehouse to a BI visualization tool, which is popularly called *Looker*.

*Looker*, a modern open-source business intelligence tool that enables users to explore and visualize data through interactive dashboards. It is used by analysts to interactively explore metrics across selected fields and across periods. It not only supports datasets to load but also supports SQL-based querying, chart creation, and filter integration, along with some metric options. The drag-and-drop option is easy to use, making it ideal for analysing large datasets quickly and securely using a web-based interface. Within the Looker web interface, a new database connection to Snowflake was created using the same credentials previously used in Airflow and dbt. This brings us one step away from drawing trends, predicting patterns, and visualising insights from the charts.

Once the connection in Looker is verified, the table *analytics.stock_summary* was added as a dataset for visualization. Using Looker's drag-and-drop chart builder, multiple insightful visualizations were created.

---

## 5. Visualization

The ETL, ELT with dbt, along with Snowflake and Looker, has brought us to the final destination where all businesses wanted to be. This is where we start drawing patterns, predicting trends, and analysing future potential. As mentioned earlier, the data pipeline focused its automation on three companies, namely, Microsoft, Apple, and Google. Let us now unfold the hidden insights for efficient business purposes.
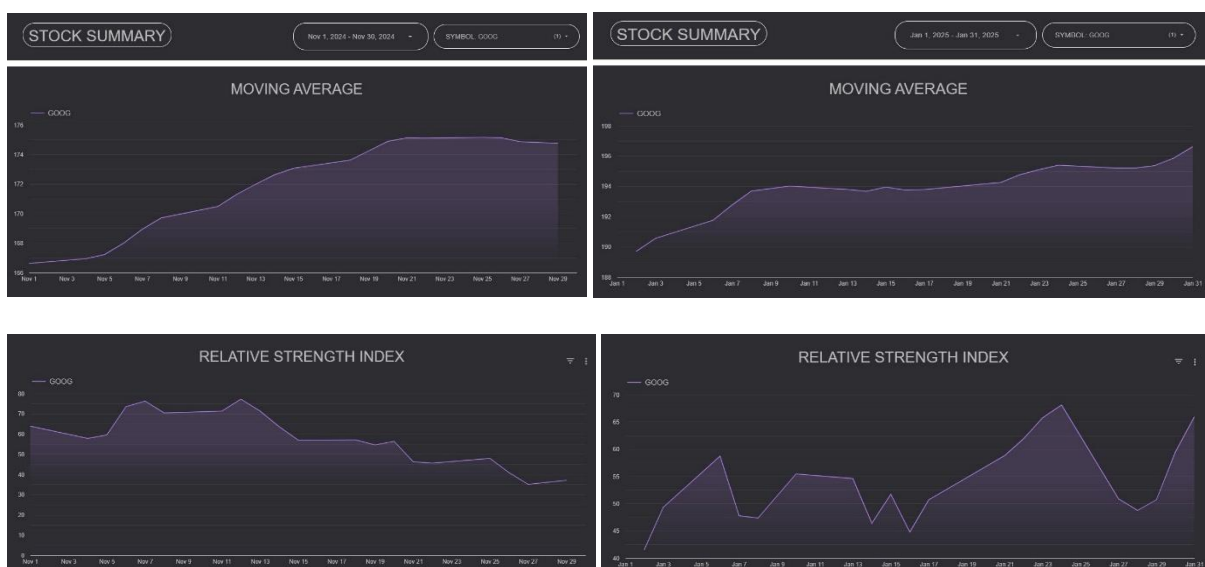
## A. Microsoft



These charts from Looker showcase Microsoft's stock performance across two time periods: **November 2024** and **January 2025**. The **Moving Average** plots indicate smoothed trends in stock prices, while November shows a gradual rise and fall, and January reflects a consistent downward trend. The **Relative Strength Index [RSI]** charts below highlight market momentum; RSI values in November fluctuated between 40-60, suggesting neutral momentum. However, in January, RSI began low and showed a gradual increase, peaking around mid-month before stabilizing. This might be due to oversold conditions early on, with recovery signs later.

These insights are especially valuable for identifying entry and exit points based on market momentum. By leveraging Looker's interactivity, users can filter by stock symbols and custom time ranges to explore deeper trends. This allows analysts to monitor changes in technical indicators like RSI in real time. Overall, the dashboard offers a streamlined way to track and interpret Microsoft's stock health visually.
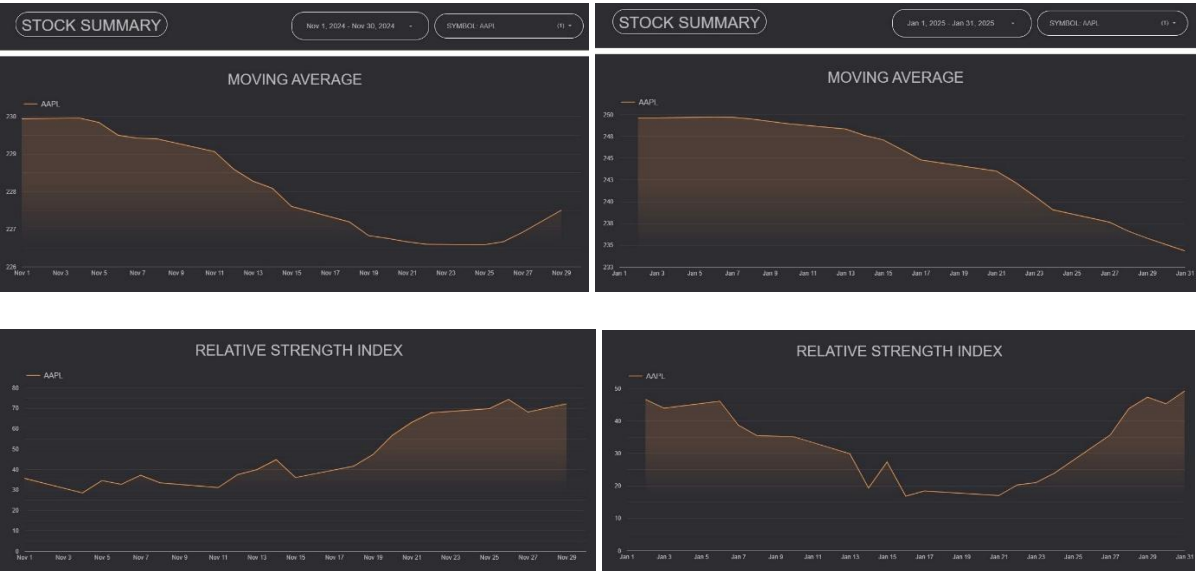
## B. Google



The visualizations present Google's stock trends over two periods: **November 2024** and **January 2025**. In November, the **Moving Average** steadily increased, showing consistent upward momentum throughout the

month. January's chart continues this trend, with a more gradual rise suggesting some market stabilization. The **Relative Strength Index (RSI)** during November started strong, hovering above 60 and even touching 75, indicating potentially overbought conditions. However, it began to decline mid-month, suggesting a slowdown in momentum. In contrast, January's RSI was more volatile, fluctuating sharply but ending on a higher note, which draws a good signal to the business industry for renewed buying interest. These dashboards make it easy to track momentum shifts and price movement patterns at a glance.

## C. Apple



These visualizations provide a clear view of Apple's stock behaviour over **November 2024** and **January 2025**. In both months, the **Moving Average** shows a downward trend, suggesting a steady decline in stock value over time. November's chart hints at a slight recovery at the end of the month, but the overall trend remained negative. January's movement was more consistent—gradually dipping throughout the month with no major upward corrections. On the **RSI [Relative Strength Index]** side, Apple started November with low momentum, but the RSI climbed steadily, even breaching 60 by the end of the month—indicating increasing buying pressure. In January, the RSI was more volatile early on, but eventually trended upward in the last 10 days, possibly hinting at a future price rebound or renewed investor interest.

---

## 6. Results

The implementation of this end-to-end data pipeline demonstrates both functional accuracy and architectural efficiency. This seamless automation helps us to draw business insights without any manual intervention at a precision rate. As the project unfolded to its final stages, the following outcomes were achieved

1. Fully automated pipeline: Extract – Transforms – Loads and Extract – Load – Transforms into Snowflake via dbt.
2. Modular DAGs: ETL and ELT split into independent Airflow DAGs.
3. Error handling: Airflow logs failures and retries as needed.
4. Clean transformations with dbt: Clear SQL logic and version control.
5. Visual insights: Traders and analysts can quickly identify stock trends, volatility, and RSI thresholds.
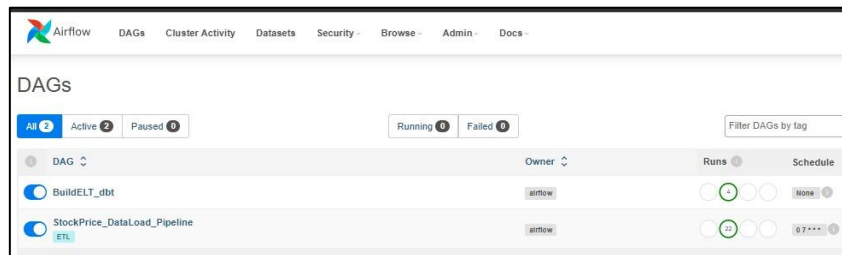
---

## 7. Discussion

The project's discussion highlights the strengths across the data pipeline components. Each tool contributes uniquely to the pipeline's functionality, reliability, and user experience. The following points summarize key observations
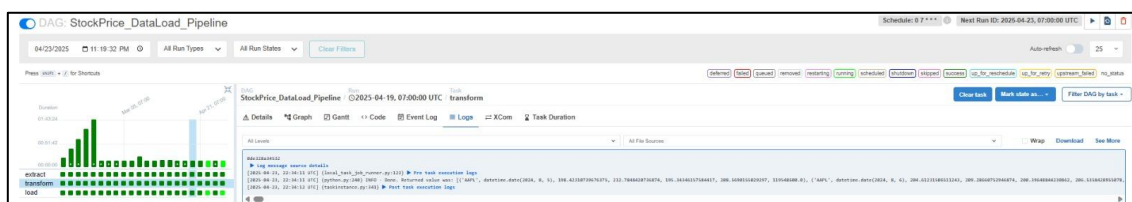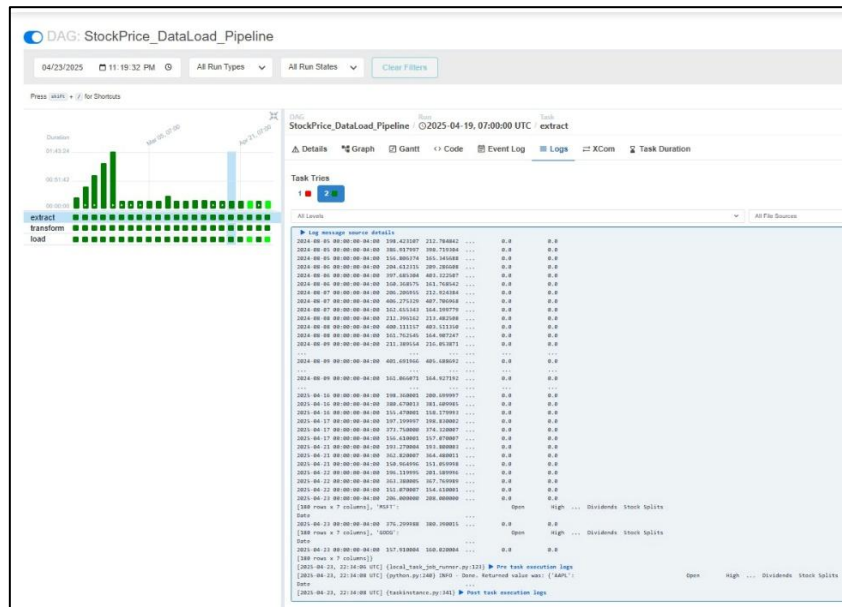
1. DBT tests are useful in verifying data quality and primary keys.
2. Airflow DAG scheduling allowed full automation without human intervention.
3. Snowflake's scalability supported fast queries even on growing datasets.
4. Looker enabled interactive exploration, but filter setup required trial and error for appropriate filters.

---

## 8. Execution Screenshots

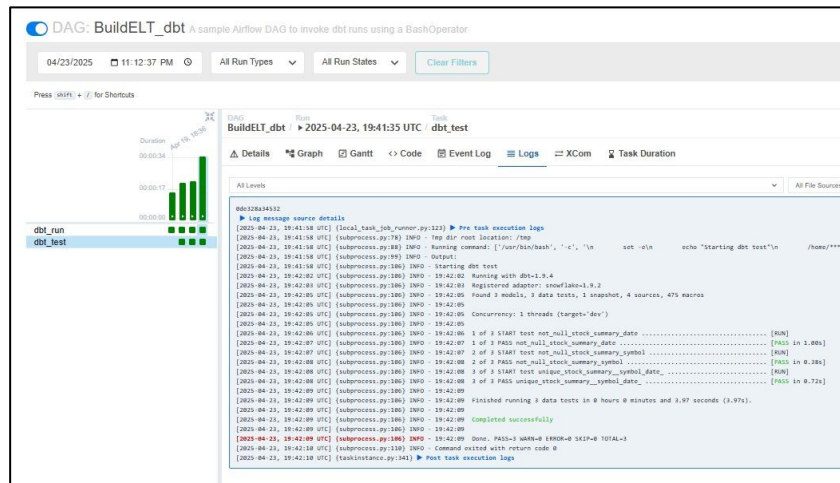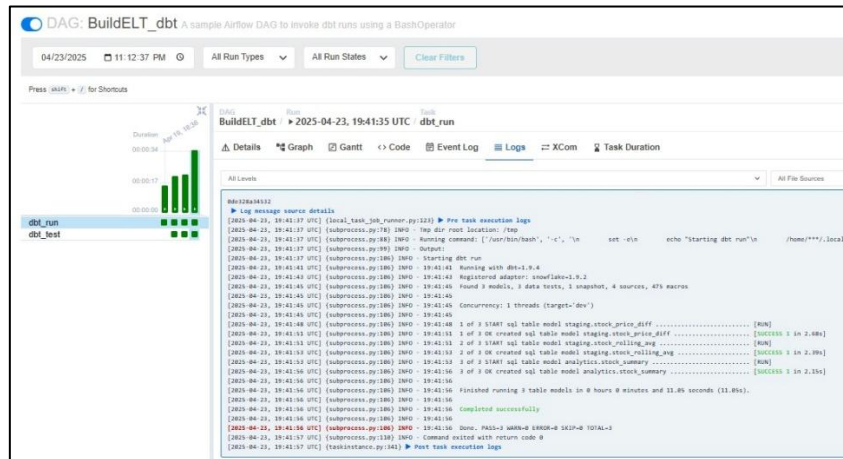### [i] Screenshot – ETL and ELT Active status



### [i] Screenshots – Extract Transform Load Airflow log screens

[i] Screenshots – ELT Airflow log screens





## 9. Conclusion

This project demonstrates how to build and deploy an end-to-end ELT pipeline using modern data engineering tools. By combining Airflow, Snowflake, dbt, and Looker, we established a scalable, automated, and insightful pipeline capable of real-time stock market analysis. Looker played a key role in transforming data into actionable insights through interactive dashboards and its unique metric feature. Its flexibility in visualizing and filtering enables analysts to explore stock trends dynamically across multiple dimensions. This modular design supports future extensions and deeper analytics.

## 10. References

[1] Apache Airflow Documentation: https://airflow.apache.org/
[2] dbt Documentation: https://docs.getdbt.com/
[3] Snowflake Documentation: https://docs.snowflake.com/
[4] Looker Documentation: https://cloud.google.com/looker/docs
[5] yFinance: https://github.com/ranaroussi/yfinance

GitHub Link: https://github.com/GoelAnshika99/Lab2-Group16/tree/main