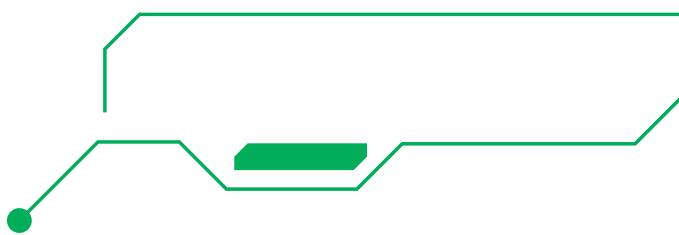


Tweet Sentiment Analysis



How do Elon Musk's tweets impact Tesla's stock price?

Presented by

Andreah Cruz,
Daniel Kim,
Jane Heng,
Lakshmi Bharathy Kumar,
Savitha Vjayarangan





PROBLEM STATEMENT

- With social media becoming a powerful communication tool, analyzing influential figures' posts can offer **insight** into how **public sentiment** shapes **financial behavior**
- Project **predicts** Tesla's next-day closing price using **engineered sentiment features** from Musk's tweets
- Dataset combines historical tweets from Elon Musk & Tesla stock price data
==> measure sentiment trends & compare them with corresponding stock movements
- This project uses:
 - Tweet data from 2010–2025 (sourced via CSV)
 - Tesla stock data pulled through Python-based ETL
 - Tools like Airflow & dbt to automate the pipeline/transform data into analytics ready formats

DATASET CONTEXT

RAW / TESLA_STOCK_DATA

Table Details Columns Data Preview Copy History Lineage

• SWAN_QUERY_WH 100 Rows • Updated 1 minute ago

	DATE	OPEN	HIGH	LOW	CLOSE	VOLUME
1	2025-05-01	280.01	290.8688	279.81	280.52	98508029
2	2025-04-30	279.9	284.45	270.78	282.16	128961057
3	2025-04-29	285.5	293.32	279.4695	292.03	108906553
4	2025-04-28	288.98	294.86	272.42	285.88	151731771
5	2025-04-25	261.69	286.85	259.63	284.95	167560688
6	2025-04-24	250.5	259.54	249.2		
7	2025-04-23	254.86	259.4499	244.43		
8	2025-04-22	230.96	242.79	229.850		
9	2025-04-21	230.26	232.21	229.850		

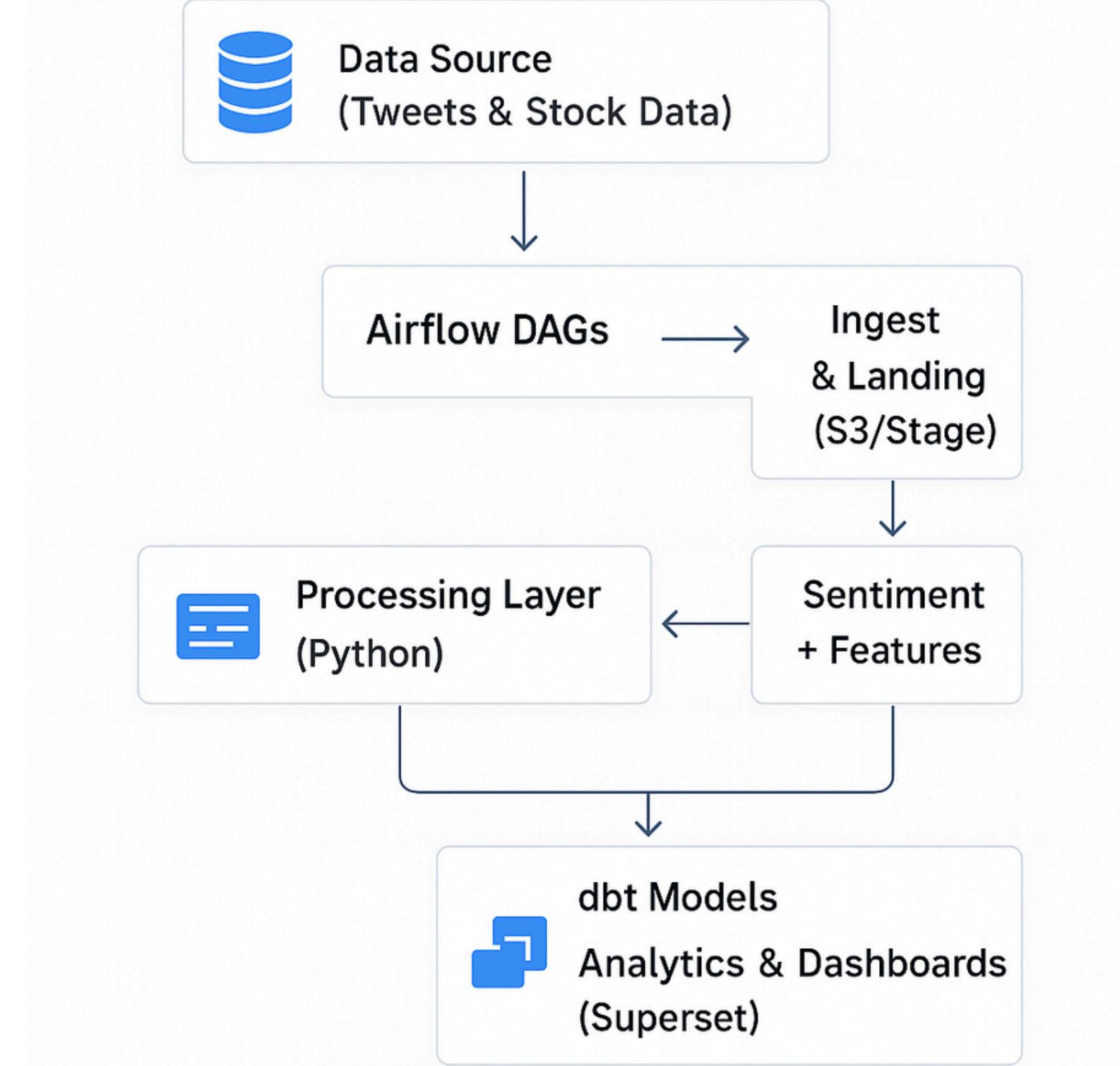
DEV / RAW / TESLA_TWEET_DATA

Table Details Columns Data Preview Copy History Lineage

• SWAN_QUERY_WH 100 of 831 Rows • Updated just now

	DATE	SENTIMENT_SCORE	WEIGHTED_SENTIMENT	TWEET_COUNT	TOTAL_LIKES
1	2023-01-01	0.1580230769	508213.3316	13	1011150
2	2023-01-02	-0.1314666667	-64236.245	6	246761
3	2023-01-03	0.07612631579	-44997.9843	19	680078
4	2023-01-04	0.2647625	93618.5409	8	966951
5	2023-01-05	0.07442727273	-99157.7027	22	1274511
6	2023-01-06	0.029625	-96557.7542	8	610642
7	2023-01-07	0.2346789474	143710.3164	19	631243
8	2023-01-08	0.1015333333	381461.3617	9	711254

- Tesla Stock Prices – Yahoo Finance API**
 - Daily open, high, low, close, volume (2010–present)
 - Collected using the **yfinance** Python package
- Elon Musk Tweets – [all_musks.csv](#) (2023-2025)**
 - Pre-collected dataset scored using the **VADER model** in Google Colab
 - Aggregated daily sentiment (sentiment score, tweet count, likes, etc.)



System ARCHITECTURE

Data Sources

Tesla stock prices & Elon Musk tweets are extracted. Sentiment is computed in Colab using VADER & stored as daily aggregates

Airflow (Ingest Layer)

Schedules and automates the daily extraction and loading of raw stock and tweet data into Snowflake staging tables

Processing Layer

Sentiment scoring was performed externally in Colab. Airflow handles scheduling; basic pre-processing logic is built into the DAG.

Snowflake Data Warehouse

Central repository for both raw and transformed tables.

dbt Transformation Layer

SQL-based models engineer features, run the regression formula, and generate predictions

Analytics & Dashboards

Superset visualizes trends, predictions, and evaluation metrics using dbt-transformed tables



View: feature_engineering_combined_features

Field Name	Data Type	Attributes	Constraints	Description
date	DATE	Not null	Primary Key with SYMBOL	Trading date
sentiment_score	FLOAT			
weighted_sentiment	FLOAT			
tweet_count	NUMBER			
total_likes	NUMBER			
lag_1_sentiment	FLOAT			
avg_3d_sentiment	FLOAT			
open_price	FLOAT			
close_price	FLOAT			
low_price	FLOAT			
high_price	FLOAT			
volume	NUMBER			
Price_change_pct	FLOAT			
high_low_range_pct	FLOAT			
Avg_7d_close	FLOAT			
weighted_tweet_impact	FLOAT			

View: evaluation_prediction_output

Field Name	Data Type	Attributes	Constraints	Description
date	DATE	Not null	Primary Key with SYMBOL	Trading date
sentiment_score	FLOAT			
weighted_sentiment	FLOAT			
tweet_count	NUMBER			
total_likes	NUMBER			
open_price	FLOAT			
close_price	FLOAT			
low_price	FLOAT			
high_price	FLOAT			
volume	NUMBER			
Price_change_pct	FLOAT			
high_low_range_pct	FLOAT			
Avg_7d_close	FLOAT			
weighted_tweet_impact	FLOAT			
Predicted_close_price	FLOAT			
Prediction_error	FLOAT			
Relative_error	FLOAT			

View: model_predict_price

Field Name	Data Type	Attributes	Constraints	Description
date	DATE	Not null	Primary Key with SYMBOL	Trading date
sentiment_score	FLOAT			
weighted_sentiment	FLOAT			
tweet_count	NUMBER			
total_likes	NUMBER			
lag_1_sentiment	FLOAT			
avg_3d_sentiment	FLOAT			
open_price	FLOAT			
close_price	FLOAT			
low_price	FLOAT			
high_price	FLOAT			
volume	NUMBER			
Price_change_pct	FLOAT			
high_low_range_pct	FLOAT			
Avg_7d_close	FLOAT			
weighted_tweet_impact	FLOAT			
Predicted_close_price	FLOAT			

DB SCHEMA

• Raw Schema

- stock_data [date, symbol] PK: daily OHLCV prices
- tweet_data [date] PK: aggregated Musk-tweet sentiment
- final_raw [date, symbol] PK: joins stock_data + tweet_data, adds price_change_pct

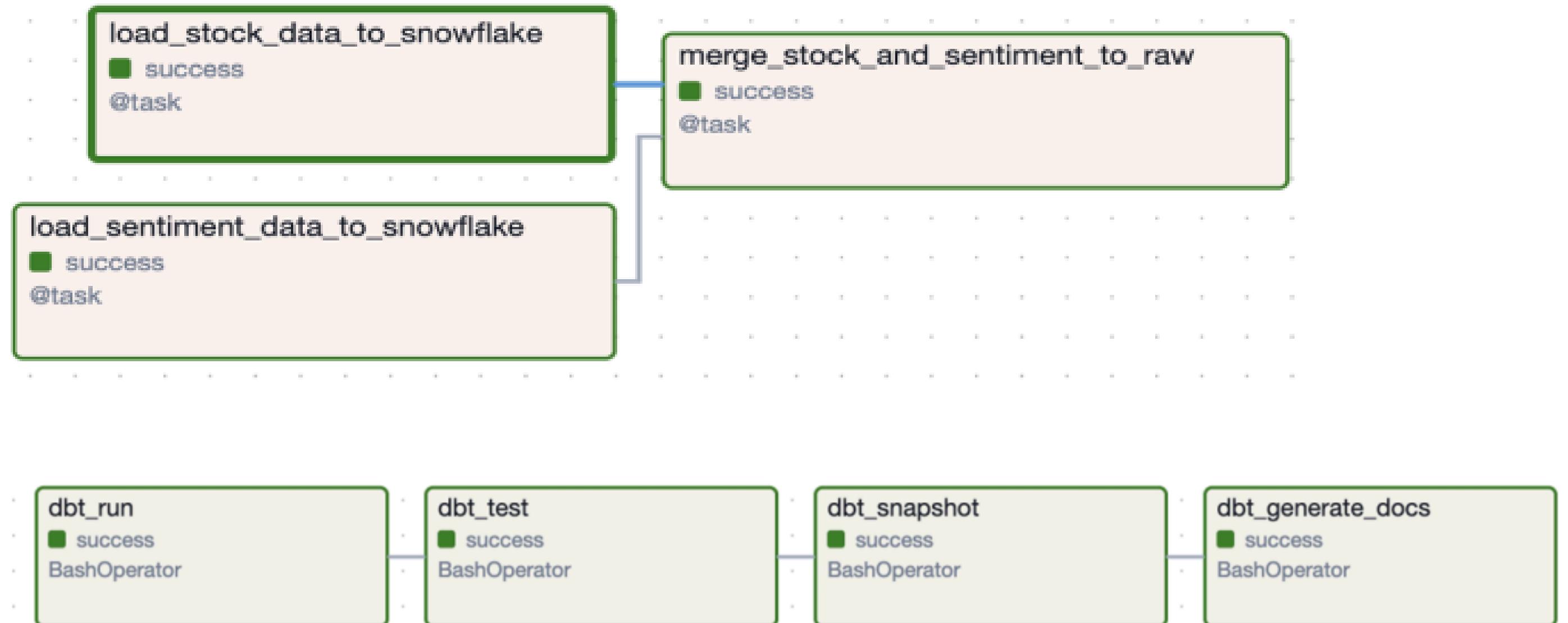
• Analytics Schema

- feature_engineering_combined_features [date, symbol] PK: rolling/lagged features
- model_predict_price [date, symbol] PK: predicted next-day close
- evaluation_prediction_output [date, symbol] PK: prediction error & relative_error

• Key Points

- Composite PK [date, symbol] ensures one row per ticker per day
- Foreign keys enforce consistency between raw tables
- Easily extensible with dim_date or dim_symbol for richer BI filters

ETL/ELT Pipeline (Airflow + dbt)

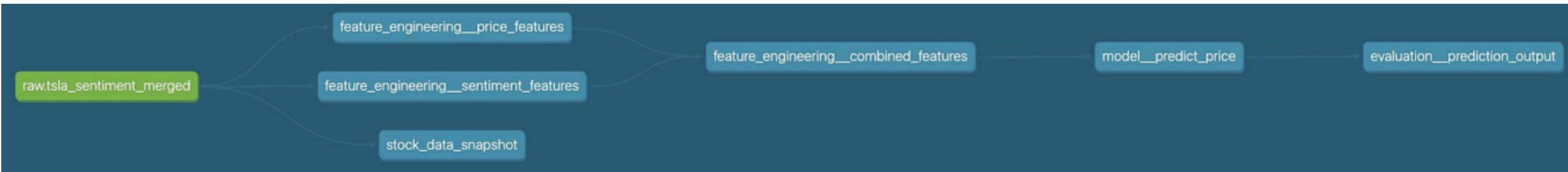
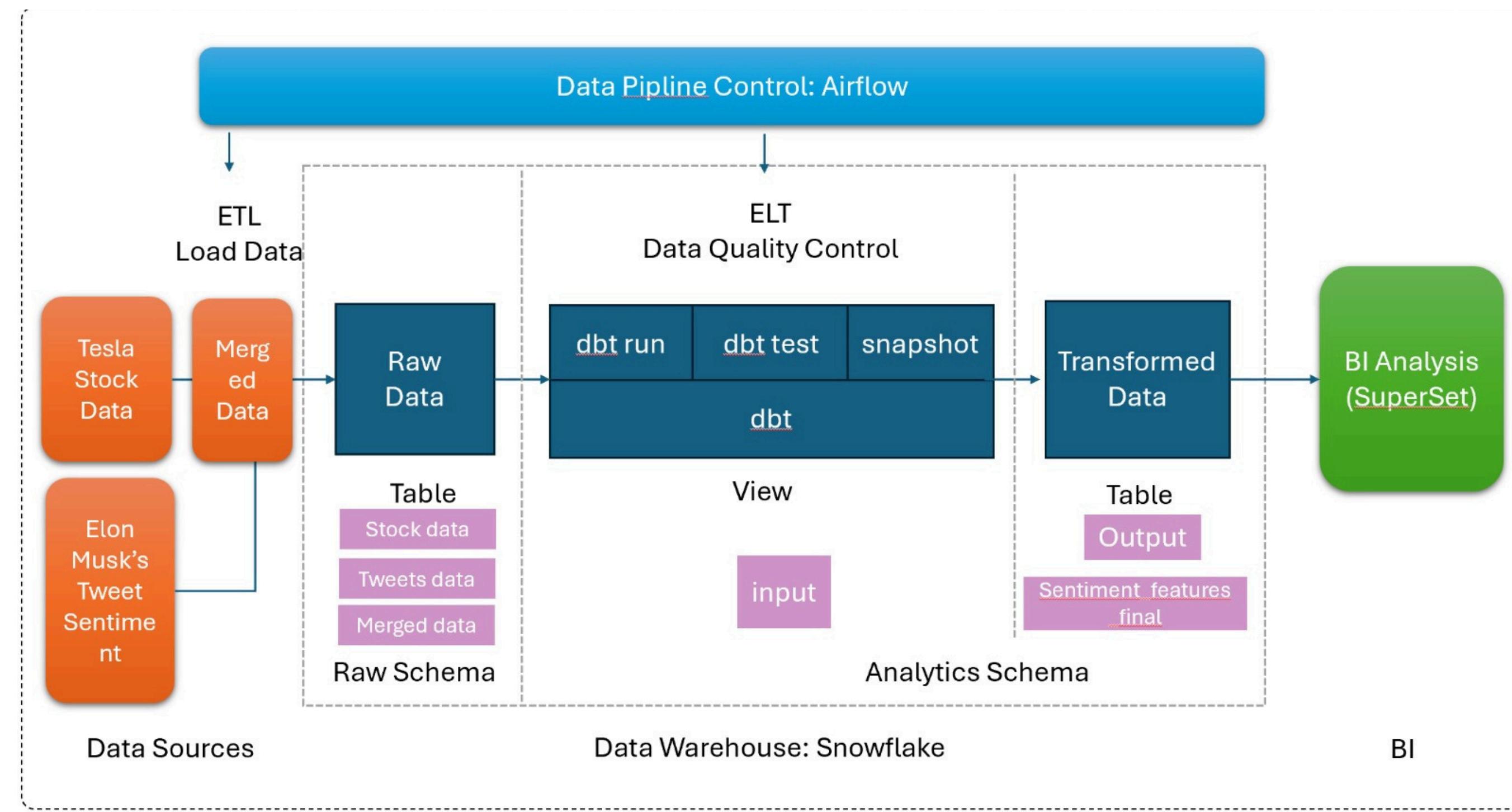


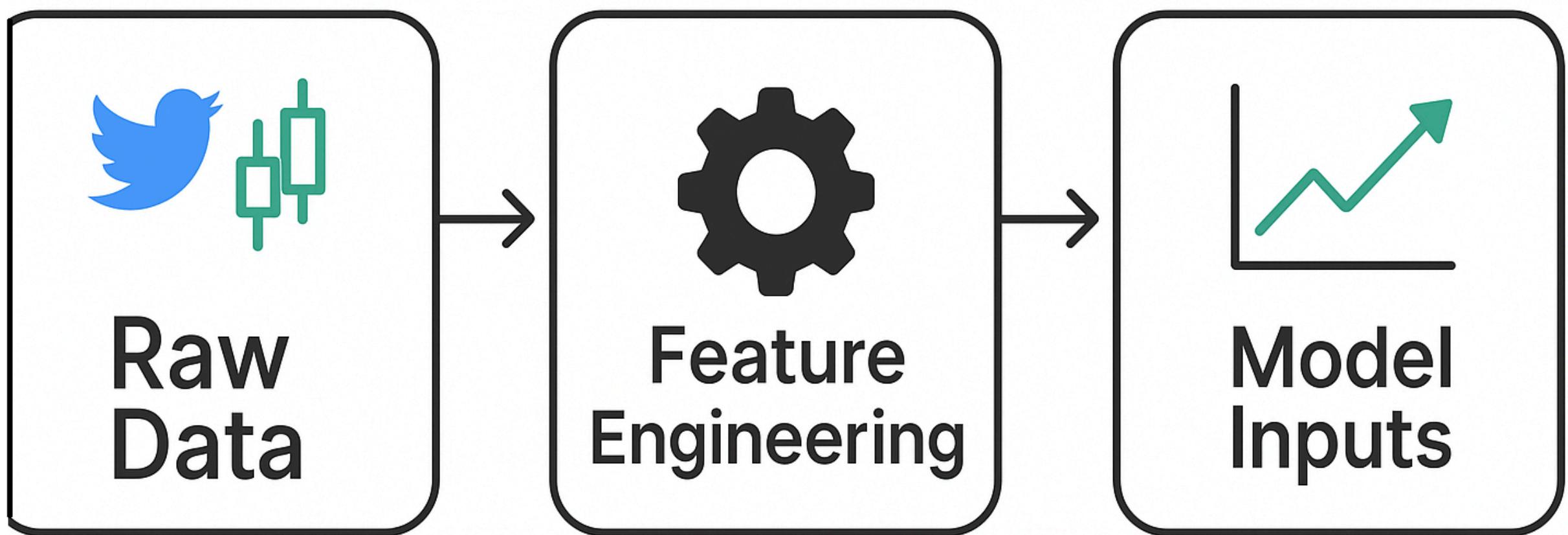
Airflow schedules:

- `extract_stock_data` – pulls latest prices
- `extract_tweet_data` – fetches Elon Musk's tweets & computes daily sentiment
- `load_stock_data_to_snowflake` – loads raw stock_data into Snowflake
- `load_sentiment_data_to_snowflake` – loads raw_tweet_data into Snowflake
- `merge_stock_and_sentiment_to_raw` – merges raw tables into final_raw_table

dbt transforms:

- `dbt run` – executes SQL models
- `dbt test` – enforces data quality checks
- `dbt snapshot` – captures historical versions of key tables for auditing





Feature Engineering

Feature Name	Coefficient (β)	Interpretation
Intercept (β_0)	-6.9047	Base value when all features are 0
SENTIMENT_S CORE	6.2242	Higher sentiment today increases price
TWEET_COUNT	0.0030	More tweets slightly increase price
LAG_1_SENTIMENT	4.9216	Yesterday's sentiment still positively influences
AVG_3D_SENTIMENT	28.5573	Recent 3-day sentiment trend has strongest impact
AVG_7D_CLOSE	1.0163	Recent price trend also contributes positively

modeling approach

VADER Sentiment Analysis Requires Python

- VADER is a Python-based NLP tool not supported in Snowflake or SQL. Using Colab allowed fast, accurate sentiment scoring on Elon Musk's tweets.

Snowflake Doesn't Support Native ML Training (for free tiers)

- Advanced ML features like FORECAST or ML_PREDICT are part of Snowflake Cortex, which requires premium licensing not available in this setup.

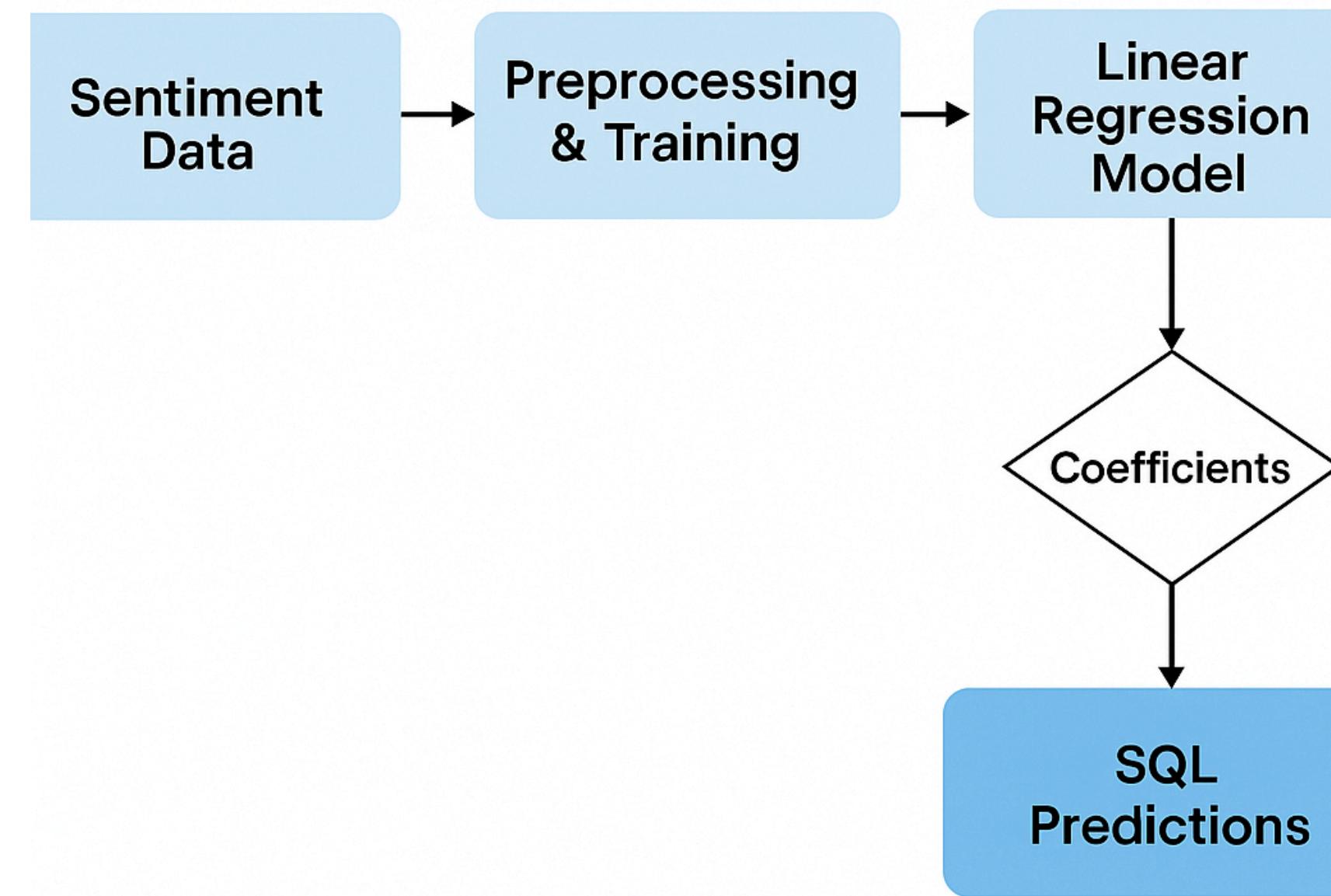
Hybrid Architecture = Best of Both Worlds

- Training and sentiment scoring were done in Python (flexible, powerful), while SQL was used for efficient, scalable in-database prediction.

In-Database Inference Enables Real-Time BI

- With predictions calculated in SQL (output.sql), results can be directly consumed by dashboards in tools like Superset without data movement.

In-Database Inference Modeling



Actual vs Predicted Line Chart

Predicted vs Actual price change



Metric	Value	Interpretation
RMSE	19.03	Avg. prediction error size in stock price units
MAE	14.58	Avg. absolute difference between predicted and actual price
R ² Score	0.9160	Model explains ~91.6% of price variation [very strong fit]

Model Evaluation



DEMO
AND
VISUALS

Deployed regression model directly within Snowflake using SQL ==> keeps entire pipeline centralized, scalable, and auditable.

- Trained in Colab, coefficients were manually imported to output.sql
- SQL calculates predicted_close_price using feature columns like:
 - lag_1_sentiment, avg_3d_sentiment, weighted_impact, etc.
- dbt model is materialized as final table used for dashboards
- Benefits:
 - No external model serving required
 - Fast, fully auditable, and easy to update
 - Leverages Snowflake's compute engine for scalability



PREDICTION EXECUTION IN SNOWFLAKE



Limitations & Challenges

- Model constraints:
 - Linear regression assumes feature independence and linearity
- Feature limitations:
 - Some engineered features (like volume) were non-informative in this setup
- Snowflake ML not used:
 - Snowflake's in-database ML was unavailable due to licensing in environment, so prediction logic was implemented manually in SQL



KEY FINDINGS:

- News or tweets sound negative ==> Tesla's stock usually dips the next day
- Positive news pushes prices up, but not by a huge amount; bigger market forces still play a role
- Price move shows up most clearly the day after the sentiment hit, so traders aren't reacting instantly

WHAT WE LEARNED/MODEL LEARNED:

- $R^2 = 0.916$ shows strong explanatory power; engineered sentiment features like `lag_1_sentiment` & `avg_3d_sentiment` were highly informative.
- Built a transparent, version-controlled pipeline using dbt & SQL—easy to scale and audit.
- Demonstrated that social sentiment can predict price, but only with thoughtful feature engineering and lag-aware modeling.



FUTURE WORK?

- Integrate intra-day prices & real-time headline streams
- Use advanced models like XGBoost or LSTM for improved accuracy
- Add more robust features: volatility, tweet engagement, macroeconomic indicators
- Expand coverage to other tickers or market sectors
- Incorporate additional sources

CONCLUSION



THANK YOU!

GROUP 9