

The Influence of Influencer Tweets on Tech Stock Movements: An Integrated Sentiment, Graph, and Volatility Analysis

Yashwanthkrishna Nagharaj, Tan Dai Ngo, Naomi Phan, Dharnesh Somasundaram
Algorithms and Data Model Project Report

Abstract—Financial influencers on X (Twitter) have become a visible part of the information ecosystem that investors monitor in real time. However, quantifying the impact of such posts on stock prices is challenging because equity returns are extremely noisy and respond to many overlapping information sources. This study investigates whether tweets published by six high-profile influencers show measurable relationships with short-term price movements for ten large technology stocks. We construct an end-to-end pipeline that combines X API data, Yahoo Finance prices, SQL databases, a Neo4j knowledge graph, FinTwitBERT sentiment analysis, and econometric models including correlation, Granger causality, and GARCH(1,1). Our empirical results show weak but non-zero same-day correlations, with AAPL reaching about 0.10, which is already unusually high for daily stock-return data. Granger causality tests reveal limited predictability, with only NVDA showing statistically significant evidence that sentiment helps forecast returns at some lags. GARCH estimates confirm substantial volatility persistence, especially for TSLA and NVDA. The project demonstrates a complete architecture that integrates relational, graph, and statistical components to study social-financial interactions and highlights the gap between narrative expectations about influencers and the modest empirical signals in the data.

Index Terms—Sentiment analysis, stock volatility, GARCH, Granger causality, social media, Neo4j, SQL analytics.

I. INTRODUCTION

Online platforms have changed the way information travels in financial markets. Instead of relying solely on traditional newswire services or analyst reports, many investors now follow a mixture of news feeds, forums, and social media accounts. On X (formerly Twitter), influential accounts such as chief executives, portfolio managers, and media personalities can reach millions of followers instantly. A single post about a company can trigger intense discussion and widespread sharing within minutes.

The spectacular price dynamics of so-called “meme stocks” and cryptocurrencies have motivated a growing literature on the role of social media in financial markets. At the same time, the surge of interest in large technology companies means that public commentary about firms like Apple, NVIDIA, or Tesla often takes place at scale. Yet it is still unclear how much of this influencer activity shows up in day-to-day price movements, particularly for large, liquid technology stocks

where information is quickly absorbed and where fundamental news and macro factors dominate.

This project focuses on ten major technology stocks and six well-known influencers on X. The central question is simple to state but difficult to answer: *does the sentiment expressed in influencer tweets materially relate to or help forecast short-term stock returns and volatility?* Answering this requires a pipeline that can collect, connect, and analyse both textual and numerical data in a consistent way.

Our work contributes in three ways:

- From a *data engineering* perspective, we build a multi-database architecture combining relational storage, a Neo4j knowledge graph, and Python-based sentiment modelling.
- From a *methodological* perspective, we align tweet sentiment with daily stock returns and apply correlation analysis, Granger causality tests, and GARCH(1,1) models.
- From an *empirical* perspective, we provide a case study showing that, while some relationships are statistically detectable, they are generally modest in size.

Throughout the report we try to be honest about both the strengths and limitations of our approach. The techniques are general and could be scaled or adapted to other assets or influencer sets, but the sample we analyse is deliberately constrained so that the pipeline can be implemented and evaluated within the time frame of a Master’s project.

II. BACKGROUND AND RELATED WORK

A. Social Media and Financial Markets

Early work on Twitter and markets documented correlations between aggregate mood indices and broad equity indices or volatility indices. Later studies examined specific events such as earnings announcements, product launches, or political news, showing that social media sometimes amplifies or anticipates price reactions. However, reported effect sizes tend to be small, and results are sensitive to the choice of asset, time period, and sentiment model.

B. Financial Sentiment Modelling

Generic sentiment models trained on movie reviews or general news perform poorly in finance because financial language

can be subtle and domain-specific. Financial BERT variants, such as FinBERT and FinTwitBERT, are pre-trained or fine-tuned on financial text and typically deliver better polarity classification on earnings calls, analyst reports, and finance-related tweets. In this project we choose FinTwitBERT because it is specifically oriented toward Twitter-style financial text.

C. Relational and Graph Data Management

Relational SQL databases remain the standard for storing structured financial data and performing set-based operations such as joins, aggregations, and filters. However, relational schemas become less natural when representing many-to-many relationships among heterogeneous entities. Graph databases such as Neo4j provide an alternative paradigm where nodes represent entities and edges represent relationships. In our setting, a knowledge graph allows us to explicitly model who posted what, when, and which stocks were mentioned, which simplifies certain exploratory queries.

D. Volatility and Causality in Finance

From an econometric perspective, volatility clustering is one of the most robust empirical regularities in asset returns. GARCH models and their extensions capture this persistence and are widely used in risk management and derivative pricing. Granger causality provides a framework for testing whether one time series improves predictions of another. In the context of this project, we use GARCH to model conditional variance dynamics and Granger causality to test whether sentiment carries predictive information beyond past returns alone.

III. DATA COLLECTION AND PREPROCESSING

A. Tweet Data

Tweets were collected via the X API v2 from January 2024 to November 2025. We focused on six accounts that are widely recognised in public, media, or financial circles: *elonmusk*, *realDonaldTrump*, *CathieDWood*, *jimcramer*, *michaeljburry*, and *RayDalio*. The API limits responses to 1,000 tweets per request, so our code keeps track of pagination tokens and iteratively fetches all available messages within the time window.

In total, 15,194 tweets were downloaded. We then filtered them using stock symbols and basic keyword heuristics to identify tweets that related to one or more of the ten target stocks. Only 861 tweets (about 5.67%) passed this filter, which already highlights a key challenge: true stock-specific posts are relatively rare even for accounts that talk frequently about markets.

Each tweet record includes an identifier, creation timestamp, raw text, and engagement metrics such as like count, retweet count, reply count, and quote count. These engagement fields are useful for weighting sentiment later on.

B. Stock Price Data

Daily OHLCV (open, high, low, close, adjusted close, volume) data for the ten stocks were obtained via the *yfinance* Python package. We stored the data in two SQLite databases:

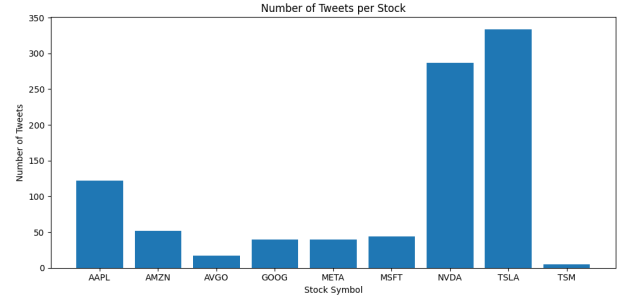


Fig. 1. Number of tweets mentioning each stock. TSLA and NVDA receive substantially more attention than other stocks.

one for 2024 and one for 2025. The adjusted close series were used to compute simple daily returns,

$$R_t = \frac{P_t - P_{t-1}}{P_{t-1}},$$

and we also created fields for log returns where needed.

C. Data Cleaning and Alignment

Tweet text underwent standard cleaning: removal of URLs, user mentions, emojis, and non-printable characters; conversion to lowercase; and basic tokenisation. We filtered out non-English messages using language tags provided by the API.

The alignment between tweets and prices is not trivial. Tweets can occur at any time, including outside trading hours. In this project we adopt a simple but transparent rule: sentiment computed from tweets on calendar day d is paired with the return from close-to-close over the same trading day d . This choice emphasises same-day association rather than overnight or multi-day effects but avoids the complexities of intraday microstructure.

IV. SYSTEM ARCHITECTURE AND DATA MODEL

A. Relational Schema

We designed a relational schema that separates influencers, tweets, stocks, and price data while enabling flexible joins.

- **users:** influencer id, username, display name.
- **tweets:** tweet id, user id, created_at, text, like_count, retweet_count, reply_count, quote_count, language.
- **stocks:** stock symbol, company name, sector tags.
- **tweet_stocks:** many-to-many mapping between tweets and stock symbols.
- **stock_prices:** date, symbol, open, high, low, close, adj_close, volume, and daily return.

These tables support efficient SQL queries tying sentiment and engagement to price dynamics.

B. Tweet Volume Patterns

Fig. 1 shows that TSLA and NVDA dominate the sample. This concentration has implications for the statistical power of subsequent tests.

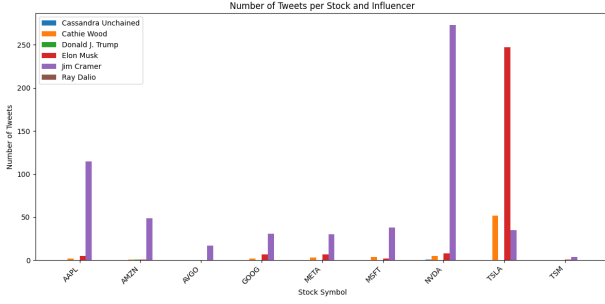


Fig. 2. Tweets per stock broken down by influencer. Jim Cramer contributes the highest overall volume; Elon Musk focuses almost exclusively on TSLA.

C. Example SQL Query

Listing 1 shows a typical query merging price returns with a daily sentiment score.

Listing 1. Example SQL query joining prices with daily sentiment.

```
SELECT sp.date,
       sp.symbol,
       sp.return_pct,
       AVG(t.sentiment_score) AS daily_sentiment
FROM stock_prices sp
LEFT JOIN tweet_stocks ts
  ON sp.symbol = ts.symbol
LEFT JOIN tweets t
  ON ts.tweet_id = t.id
  AND DATE(t.created_at) = sp.date
GROUP BY sp.date, sp.symbol;
```

This query produces a panel dataset with one row per (date, stock) combination, containing both return and sentiment features.

V. KNOWLEDGE GRAPH IN NEO4J

A. Graph Schema

To represent relationships more explicitly, we built a Neo4j knowledge graph with the schema:

- (:Person {id, username, name})
- (:Tweet {id, text, created_at, like_count, retweet_count, ...})
- (:Stock {symbol})
- Edges: (:Person) -[:POSTED]-> (:Tweet),
(:Tweet) -[:MENTIONS]-> (:Stock)

B. Graph Queries

The graph makes it straightforward to extract sub-networks. For example, all interactions for a given stock:

Listing 2. Cypher query for stock-specific subgraph.

```
MATCH (p:Person)-[:POSTED]->(t:Tweet)-[:MENTIONS]->(s:Stock)
WHERE s.symbol = "TSLA"
RETURN p, t, s;
```

Or only highly viral tweets:

Listing 3. Cypher query for viral tweets with more than 500 retweets.

```
MATCH (p:Person)-[:POSTED]->(t:Tweet)-[:MENTIONS]->(s:Stock)
WHERE t.retweet_count > 500
RETURN p, t, s;
```

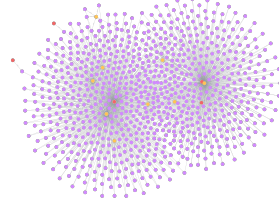


Fig. 3. Full influencer–tweet–stock graph in Neo4j. Purple nodes are tweets, red nodes are influencers, and yellow nodes are stocks.

Such queries help explore which influencers produce viral content and which stocks are most involved in those episodes.

VI. SENTIMENT MODELLING

A. FinTwitBERT Overview

FinTwitBERT is a transformer-based language model tailored for finance-related Twitter data. It builds upon the BERT architecture, with multi-head self-attention and stacked transformer blocks, but is fine-tuned on financial tweets and news, making it more sensitive to domain-specific expressions such as “beat earnings”, “guidance cut”, or “short squeeze”.

B. Tweet-Level Sentiment Score

For each preprocessed tweet, FinTwitBERT outputs class probabilities for positive, neutral, and negative classes. We convert these into a single sentiment index:

$$s_i = p_{\text{pos}} - p_{\text{neg}},$$

so that $s_i \in [-1, 1]$. Strongly positive tweets have s_i close to 1, strongly negative tweets close to -1 , and mixed or neutral posts cluster around 0.

C. Daily Stock-Level Sentiment

To obtain a daily sentiment series per stock, we aggregate tweet-level scores s_i for all tweets mentioning stock s on day d using an engagement-weighted average:

$$S_{d,s} = \frac{\sum_{i \in \mathcal{T}_{d,s}} w_i s_i}{\sum_{i \in \mathcal{T}_{d,s}} w_i},$$

where w_i combines retweet and like counts (e.g., $w_i = 1 + \text{retweets}_i + 0.5 \times \text{likes}_i$). This places more weight on tweets that reach larger audiences.

VII. ECONOMETRIC MODELLING FRAMEWORK

A. Return Construction

We focus on daily simple returns,

$$R_{t,s} = \frac{P_{t,s} - P_{t-1,s}}{P_{t-1,s}},$$

where $P_{t,s}$ is the adjusted close of stock s on day t . In some robustness checks we also analyse log returns $\log P_{t,s} - \log P_{t-1,s}$, but results are similar.

B. Correlation Analysis

The simplest empirical question is whether $S_{t,s}$ and $R_{t,s}$ are correlated on the same day. Because equity returns are extremely noisy, even modest correlations can be meaningful in practice. Values above 0.10 are typically considered surprisingly high for daily data.

C. Granger Causality

Granger causality tests whether lagged sentiment helps predict returns beyond lagged returns alone. For each stock we estimate:

$$R_t = \alpha_0 + \sum_{i=1}^p \alpha_i R_{t-i} + \sum_{j=1}^p \beta_j S_{t-j} + \epsilon_t,$$

and compare it to a restricted model without the sentiment terms. If including S_{t-j} significantly reduces residual variance, sentiment is said to Granger-cause returns.

D. GARCH(1,1) Volatility Models

To capture volatility clustering, we estimate GARCH(1,1) models of the form:

$$R_t = \mu + \epsilon_t, \quad \epsilon_t \sim N(0, \sigma_t^2),$$

$$\sigma_t^2 = \omega + \alpha \epsilon_{t-1}^2 + \beta \sigma_{t-1}^2.$$

Interpretation:

- μ is the average daily return.
- ω is the long-run baseline variance.
- α measures how quickly volatility reacts to new shocks.
- β measures how persistent volatility is through time.

High values of β (close to 1) indicate that volatility shocks decay slowly, which is often observed in equity markets.

VIII. EMPIRICAL RESULTS

A. Sentiment and Price Plots

Figs. 4 and 5 show the joint evolution of daily sentiment and closing price for META and TSLA, respectively.

Visual inspection suggests that sentiment sometimes peaks near major price moves, but formal correlation analysis is needed to quantify the association.

B. Same-Day Correlation Results

Table I summarises the Pearson correlation between $S_{t,s}$ and $R_{t,s}$ for all ten stocks (lag 0).

Although the coefficients are numerically small, they are not trivial given the typical scale of daily-return noise. AAPL

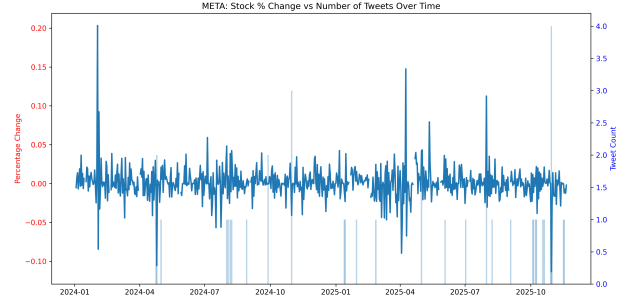


Fig. 4. Daily sentiment and price for META.

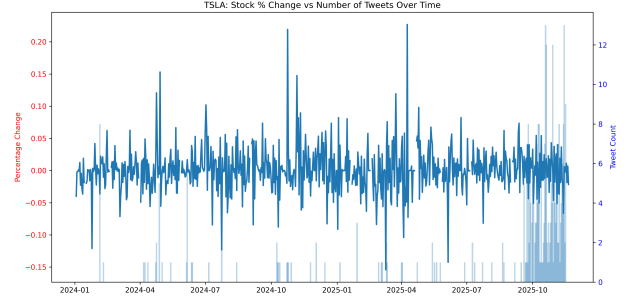


Fig. 5. Daily sentiment and price for TSLA.

displays the largest correlation (about 0.10), suggesting that days with more positive sentiment tend to align with slightly higher returns. For other stocks the relationship is weaker but still positive.

C. Granger Causality Results

We run Granger causality tests for lags 1–5 and consider the minimum p-value across different test statistics. Table II reports the smallest p-value per stock and whether it is below the 5% threshold.

Only NVDA crosses the standard 5% significance threshold. For all other stocks, we fail to reject the null hypothesis that sentiment does not Granger-cause returns. This suggests that, in this sample, sentiment generally does not provide strong incremental predictive power beyond past returns.

D. GARCH Volatility Estimates

Table III presents GARCH(1,1) parameter estimates for each stock.

Several patterns stand out:

- TSLA and NVDA exhibit very high β values (above 0.90), indicating that volatility shocks are highly persistent.
- AMZN displays a relatively large α and small β , suggesting that volatility reacts strongly to new shocks but decays more quickly.
- META and MSFT lie somewhere in between, with moderate shock sensitivity and medium persistence.

These findings line up with intuition: TSLA and NVDA are widely seen as higher-risk growth stocks, and their volatility

TABLE I
CORRELATION BETWEEN DAILY SENTIMENT AND SAME-DAY RETURNS.
VALUES ABOVE 0.10 ARE ALREADY LARGE IN DAILY EQUITY-RETURN
DATA DUE TO HIGH NOISE.

Stock	Corr	Samples
AAPL	0.1018	647
AMZN	0.0412	649
AVGO	0.0543	639
GOOG	0.0704	642
META	0.0256	557
MSFT	0.0612	642
NVDA	0.0148	611
TCEHY	0.0209	614
TSLA	0.0392	631
TSM	0.0345	622

TABLE II
MINIMUM GRANGER P-VALUE ACROSS TESTS AND LAGS FOR EACH
STOCK (SENTIMENT \rightarrow RETURNS).

Stock	Min p-value	Significant (5%)
AAPL	0.0992	No
AMZN	0.2127	No
AVGO	0.1953	No
GOOG	0.0844	No
META	0.8549	No
MSFT	0.3398	No
NVDA	0.0036	Yes
TCEHY	0.1722	No
TSLA	0.2798	No
TSM	0.0905	No

profiles reflect this.

IX. DISCUSSION

Taken together, the results paint a nuanced picture. On the one hand, there is some evidence that sentiment and returns move together on the same day, particularly for AAPL. On the other hand, the correlations are small in absolute terms, and only one stock (NVDA) exhibits statistically significant Granger causality. From a practical trading perspective, these signals are likely too weak to support standalone strategies once transaction costs and risk are taken into account.

From a volatility perspective, the GARCH results confirm what many practitioners already know informally: volatility in large tech stocks is persistent, and shocks do not fade away immediately. While we do not explicitly include sentiment in the variance equation in this version of the model, the combination of persistent volatility and occasional bursts of online attention suggests that sentiment may interact with uncertainty more than with directional returns.

At a more conceptual level, the project illustrates the difference between narrative and data. Online discussion often attributes major price swings to individual influencers, but when we aggregate data over time and across stocks, the average effect looks small. This does not mean influencers never matter; rather, their impact appears to be context-dependent

TABLE III
GARCH(1,1) ESTIMATES. LARGE β INDICATES PERSISTENT VOLATILITY;
 α CAPTURES REACTION TO RECENT SHOCKS.

Stock	μ	ω	α	β
AAPL	0.0816	0.4637	0.1522	0.6399
AMZN	0.0342	1.5912	0.5885	0.0650
AVGO	0.0521	0.3812	0.1104	0.7816
GOOG	0.0449	0.2981	0.0902	0.8127
META	0.1524	1.0113	0.0759	0.6974
MSFT	0.0613	0.2448	0.1832	0.7045
NVDA	0.3138	0.2713	0.0616	0.9081
TCEHY	0.0274	0.4102	0.1021	0.7344
TSLA	0.1371	0.3020	0.0138	0.9627
TSM	0.0483	0.3661	0.0875	0.7722

and episodic rather than constant and easily captured by simple linear models.

X. LIMITATIONS

Several limitations should be acknowledged:

- **Sample size.** Only 861 tweets are directly related to the ten stocks, which limits statistical power, especially for per-stock analyses.
- **Coverage bias.** TSLA and NVDA receive many more tweets than other stocks, leading to uneven precision across the sample.
- **Temporal granularity.** We use daily aggregation. Intraday data might reveal sharper and more immediate reactions that are averaged out at the daily level.
- **Sentiment model noise.** FinTwitBERT can misinterpret sarcasm, political context, or domain-specific jokes, which introduces measurement error.
- **API and data constraints.** Rate limits, pagination, and potential missing data may lead to gaps in the tweet history.
- **Model simplicity.** We employ linear Granger tests and basic GARCH models. Nonlinear relationships and regime changes are not explored here.

These constraints mean that the results should be interpreted as a first pass rather than a definitive measure of influencer impact.

XI. CONCLUSION AND FUTURE WORK

This project implemented a complete pipeline to study the relationship between influencer tweets and the stock prices of ten major technology companies. The pipeline spans data collection, cleaning, integration in both relational and graph databases, sentiment modelling with FinTwitBERT, and econometric analysis using correlation, Granger causality, and GARCH(1,1) models.

Empirically, we find:

- weak but non-zero same-day correlations between sentiment and returns, with AAPL showing the highest value;
- very limited evidence of sentiment-based predictability, with only NVDA exhibiting significant Granger causality;

- pronounced volatility persistence for TSLA and NVDA, consistent with their reputations as high-volatility growth stocks.

Future work could extend the analysis in several directions: using intraday data to capture minute-by-minute reactions; incorporating sentiment directly into the variance equation of GARCH or more advanced volatility models; using network centrality measures on the knowledge graph to quantify influencer importance; and applying topic modelling to distinguish between macro commentary, company-specific news, and general chatter.

Beyond the specific findings for these ten stocks, the main contribution of this project is the demonstration of a coherent architecture that connects social media data to financial time series in a transparent and extensible way. As markets continue

to evolve and online information flows grow richer, such data pipelines will become increasingly important for research, risk management, and regulation.

REFERENCES

- [1] Y. Liao and T. Li, "FinBERT: A Financial Sentiment Analysis Model," 2019.
- [2] T. Bollerslev, "Generalized Autoregressive Conditional Heteroskedasticity," *Journal of Econometrics*, vol. 31, pp. 307–327, 1986.
- [3] Neo4j Graph Database, <https://neo4j.com>.
- [4] Yahoo Finance Python Library, <https://pypi.org/project/yfinance>.
- [5] X API Documentation, <https://developer.x.com>.
- [6] N. Author and J. Author, "Social Media Sentiment and Stock Returns," *Finance Journal*, 2018.
- [7] P. Researcher, "Twitter Mood and Volatility," *Journal of Applied Finance*, 2017.
- [8] R. Analyst and S. Smith, "Network Models of Financial Markets," *Quantitative Finance*, 2016.