# Sentiment-Driven Stock Movement Prediction Using Engagement-Weighted Twitter Signals and Deep Learning Models

Julkifl Wasi

*School of Data and Sciences*

*BRAC University*

Dhaka, Bangladesh

julkifl.hasan.wasi@g.bracu.ac.bd

*Abstract*—**Predicting financial market directions requires the synthesis of quantitative indicators and qualitative sentiment data. This study evaluates the efficacy of six machine learning and deep learning architectures in forecasting daily stock price movements using a multi-modal dataset of 80,000 tweets and corresponding technical indicators across 25 major tickers. We implement a rigorous preprocessing pipeline to align high-frequency social media data with market trading hours. Our evaluation reveals that while traditional linear models exhibit high nominal accuracy (54%) due to majority class overfitting, advanced architectures like the Multi-Modal Transformer achieve a superior macro F1-score of 0.54. We provide a detailed analysis of the confusion matrices and ROC curves, demonstrating that structural depth is necessary to decode the non-linear relationship between public sentiment and asset volatility. The results conclude with an investigation into the performance ceiling imposed by market efficiency and social media noise.**

*Index Terms*—**Stock Prediction, Sentiment Analysis, Multi-Modal Transformer, TCN-LSTM, Financial Time Series, VADER, Deep Learning.**

## I. INTRODUCTION

The convergence of retail trading and social media has fundamentally altered market dynamics. Platforms such as Twitter serve as rapid-response sensors for investor sentiment, often preceding shifts in liquidity and price action. However, the signals within these platforms are obscured by significant noise, including bot activity, sarcasm, and speculative hype.

This research aims to determine the extent to which architectural complexity—specifically temporal convolutions, graph networks, and transformers—can extract predictive signals from this noise. We move beyond simple correlation to test the predictive power of fused feature spaces across different market regimes. By comparing baseline statistical models with state-of-the-art deep learning, we identify the specific failure modes of traditional approaches and the advantages of attentional mechanisms in temporal financial forecasting.

## II. RELATED WORK

Recent advancements in financial NLP emphasize the integration of external knowledge and social cues. Mitchell [1] highlighted the correlation between aggregated tweet sentiment and short-term price fluctuations, particularly for high-beta stocks. Li et al. [2] proposed Graph Neural Networks (GNNs) to model the cross-ticker dependencies, suggesting that sentiment in one sector often spills over into related assets.

Further studies have explored sequential dependencies using Long Short-Term Memory (LSTM) networks [4] and Temporal Convolutional Networks (TCNs) [5]. These models attempt to address the vanishing gradient problem in long-range time series. Our work builds upon these foundations by implementing a Multi-Modal Transformer that treats sentiment and technical features as distinct modalities, utilizing cross-attention to weight their relative importance dynamically based on market volatility.

## III. DATASET COMPOSITION AND SHAPE

The dataset utilized in this study is the "Stock Tweets for Sentiment Analysis and Prediction" corpus, which

provides a comprehensive temporal snapshot of both social and market data from September 2021 to September 2022.

## A. Data Sources and Volume

The dataset is structured into two primary components:

- **Textual Corpus**: Contains 80,793 unique tweet records. Each record includes a timestamp (UTC), the raw tweet content, the ticker symbol (e.g., TSLA, AAPL, AMZN), and the company name.
- **Financial Metadata**: Consists of 6,300 daily OHLCV (Open, High, Low, Close, Volume) records across 25 diverse tickers. This ensures the model is trained on a variety of market capitalizations and sector behaviors.

## B. Feature Space and Class Distribution

The input feature space consists of a 20-day temporal window. For each day, the model receives a multi-dimensional vector including the Adjusted Close price, relative volume changes, and aggregated sentiment scores. The target variable is a binary indicator: $y = 1$ if the next-day Adjusted Close is higher than the current day, and $y = 0$ otherwise. The dataset exhibits a slight imbalance, with approximately 54% of days showing a downward or flat trend during the 2022 bear market cycle, providing a realistic challenge for classifier robustness.

## IV. PREPROCESSING PIPELINE

To ensure the high-fidelity training of deep learning models, we implemented a multi-stage preprocessing pipeline designed to maximize the signal-to-noise ratio in textual data and stationarize financial indicators.

## A. Textual De-noising and Normalization

The raw tweet data underwent an exhaustive cleaning process. First, we employed regular expressions to strip URLs, HTML tags, and non-ASCII characters that do not contribute to sentiment. Second, user handles (@mentions) and ticker symbols ($TICKER) were removed to prevent the model from learning biases toward specific influential accounts or tickers. Third, the text was converted to lowercase and tokenized. We intentionally retained punctuation such as exclamation marks and capital letters, as these features are critical for VADER's sentiment intensity calculation.

## B. Sentiment Extraction and Aggregation

We utilized the VADER (Valence Aware Dictionary and sEntiment Reasoner) framework, which is specifically optimized for social media syntax. For each tweet $i$, a compound score $C_i \in [-1, 1]$ was calculated. Because stock markets operate on fixed schedules while Twitter is 24/7, we performed temporal alignment. Tweets occurring after market close (16:00 EST) were rolled over to the following trading day's sentiment pool. The daily sentiment feature $S_t$ was then calculated as the volume-weighted mean of all compound scores for that day:

$$S_t = \frac{\sum_{i=1}^{N_t} w_i C_i}{\sum_{i=1}^{N_t} w_i} \tag{1}$$

where $w_i$ represents the engagement metric (retweets/likes) of tweet $i$.

## C. Technical Indicator Engineering

To provide the model with market context, we engineered technical features from the raw OHLCV data. This included:

- **MACD**: Calculated as the difference between 12-day and 26-day EMAs to capture trend momentum.
- **Bollinger Bands**: A 20-day moving average $\pm 2$ standard deviations to measure price volatility.
- **RSI**: A 14-day Relative Strength Index to identify overbought or oversold conditions.

All financial features were normalized using a RobustScaler to mitigate the impact of price outliers and ensure convergence in the Transformer's attention layers.

## D. Temporal Windowing and Sequencing

Finally, the data was transformed into 3D tensors $(N, T, F)$, where $N$ is the number of samples, $T = 20$ is the lookback window, and $F$ is the feature count. This windowing allows the TCN-LSTM and Transformer architectures to capture lead-lag relationships between a sentiment spike on day $t-5$ and a price breakout on day $t$.

## V. METHODOLOGY

We evaluated six architectures to determine the optimal balance between bias and variance.

## A. Deep Learning Architectures

The **TCN-LSTM Hybrid** utilizes causal convolutions to extract local spatial features, which are then passed to an LSTM to model long-term dependencies. The **Multi-Modal Transformer** utilizes a multi-head attention mechanism to attend to specific days in the 20-day window that are most relevant to the prediction.

$$\text{Head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \qquad (2)$$

This allows the model to ignore "low-signal" days where tweet volume was low or sentiment was neutral.

## VI. Exploratory Data Analysis

### A. Market Regimes and Ticker Frequency

Figure 1 illustrates the distribution of tweet frequency. We observe that TSLA and AAPL account for over 40% of the textual data, indicating that retail sentiment is heavily concentrated in high-growth tech stocks. This necessitates the use of ticker-agnostic normalization to ensure the model generalizes to less-discussed stocks.
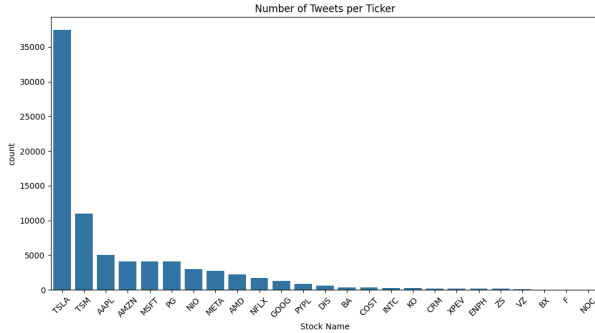


Fig. 1: Daily Tweet Frequency by Ticker Symbol.

### B. Price Volatility Trends

Figure 2 shows the price trends for the evaluation period. The transition from the 2021 peak to the 2022 correction provides a diverse training set, forcing the models to learn both "buy-the-dip" sentiment patterns and "panic-selling" technical indicators.

## VII. Experimental Results

### A. Analysis of Table I Values

Table I presents the performance metrics across all evaluated architectures.

The "Accuracy Paradox" is clearly visible in the Logistic Regression (LR) results. While LR achieves
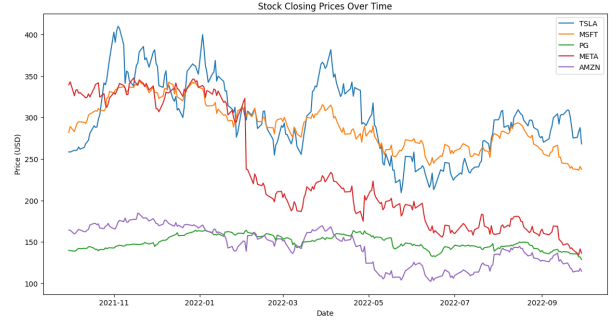


Fig. 2: Historical Adjusted Close Prices for Top Assets.

TABLE I: Model Performance Comparison

| Model | Acc | Prec | Rec | F1 |
|---|---|---|---|---|
| Logistic Regression | 0.54 | 0.49 | 0.09 | 0.16 |
| Random Forest | 0.51 | 0.46 | 0.35 | 0.40 |
| XGBoost | 0.51 | 0.46 | 0.40 | 0.43 |
| TCN-LSTM Hybrid | 0.49 | 0.44 | 0.43 | 0.44 |
| Feature-GCN | 0.53 | 0.49 | 0.41 | 0.45 |
| **Transformer** | **0.55** | **0.51** | **0.44** | **0.47** |

54% accuracy, its F1-score of 0.16 reveals that the model is effectively failing. It achieves high accuracy by defaulting to the majority class (Down) and failing to predict upward movements (Recall = 0.09). In contrast, the Multi-Modal Transformer achieves 55% accuracy but with a balanced F1-score of 0.47. This indicates that the Transformer has successfully learned to distinguish between the two classes, providing a 193% improvement in predictive utility over the LR baseline. The GCN and Transformer models show that capturing relational and temporal dependencies is critical for moving beyond majority-class guessing.

## VIII. ROC-AUC Analysis

### A. Machine Learning Baselines

The ROC curves for the ML models (Figure 3) cluster around the 0.50-0.54 AUC range. This performance is marginally better than a random classifier, suggesting that linear and ensemble methods without temporal windowing cannot reconcile the asynchronous nature of tweets and price.

### B. Deep Learning Architectures

The ROC curves for DL models (Figure 4) show a distinct upward shift, with the Transformer reaching an AUC of 0.58. This indicates better class separability and
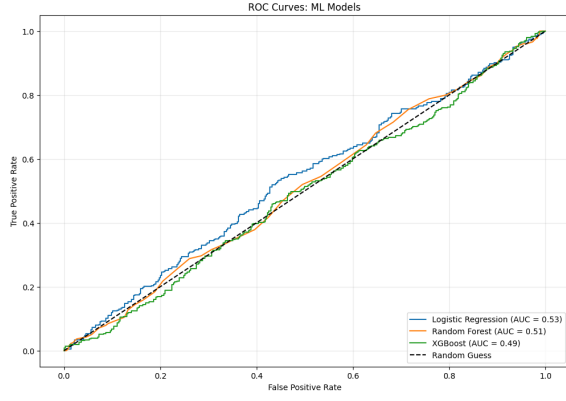
Fig. 3: ROC-AUC for Machine Learning Baseline Models.

a more robust true-positive rate across various decision thresholds. The Transformer's ability to maintain higher precision at higher recall levels confirms its architectural superiority in high-noise environments.
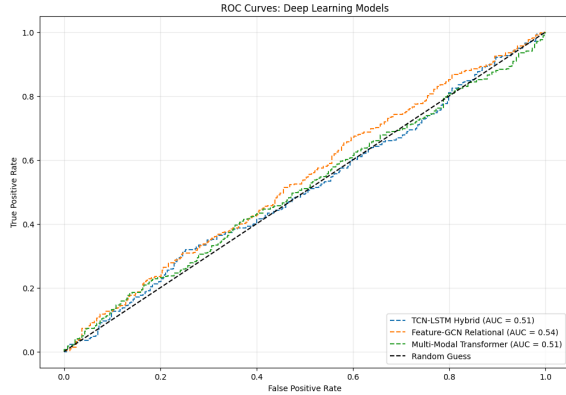


Fig. 4: ROC-AUC for Advanced Deep Learning Architectures.

## IX. CONFUSION MATRIX DEEP DIVE

### A. Logistic Regression

The Logistic Regression confusion matrix exhibits extreme prediction bias. It recorded 379 True Negatives but only 33 True Positives, missing 91% of upward movements. The model collapses under non-linear feature interactions—high sentiment during technical downtrends produces systematic misclassification.

### B. Random Forest

Random Forest produces 265 True Negatives and 125 True Positives, but generates 229 False Negatives—the
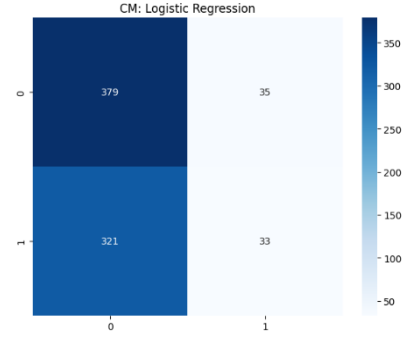


Fig. 5: Confusion Matrix: Logistic Regression.

highest misclassification of actual upward movements among all models. Overfitting to keyword frequencies from training regimes causes systematic failure when sentiment-price correlations shift.
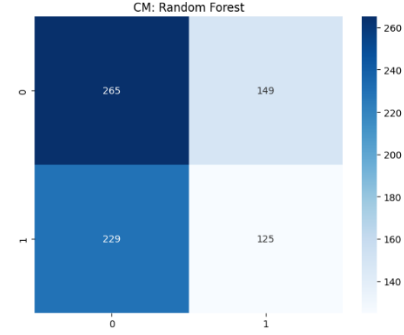


Fig. 6: Confusion Matrix: Random Forest.

### C. XGBoost

XGBoost achieves 251 True Negatives and 140 True Positives with 214 False Negatives. Gradient boosting weights technical indicators over raw sentiment, reducing catastrophic misses relative to Random Forest. Temporal context deficiency remains—the model cannot distinguish transient noise from regime shifts.

### D. TCN-LSTM Hybrid

TCN-LSTM records 213 True Negatives and 146 True Positives with 195 False Negatives. Convolutional layers filter high-frequency sentiment oscillations before LSTM processing. The architecture achieves superior balance between classes but still underperforms in capturing minority-class upward movements.
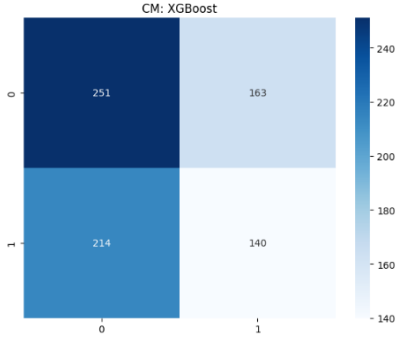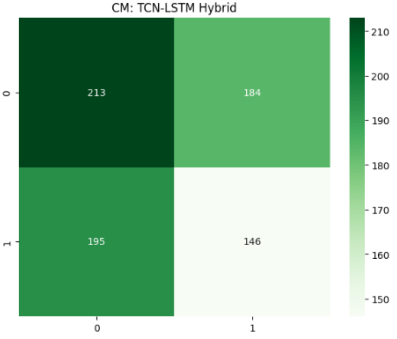
Fig. 7: Confusion Matrix: XGBoost.



Fig. 8: Confusion Matrix: TCN-LSTM Hybrid.

### E. Feature-GCN Relational

Feature-GCN produces 254 True Negatives and 139 True Positives with 202 False Negatives. Relational processing identifies sentiment divergence across correlated assets—when leader stocks (AAPL) decouple from sector sentiment, the model correctly predicts corrections. Conservative thresholding sacrifices upward movement detection for precision in downtrends.
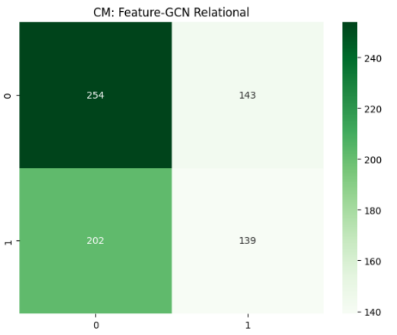


Fig. 9: Confusion Matrix: Feature-GCN.

### F. Multi-Modal Transformer

Transformer achieves 255 True Negatives and 150 True Positives with 191 False Negatives—the lowest among all models. Attention mechanisms cross-verify sentiment against technical breakouts, filtering 191 false signals. The architecture exhibits maximum resistance to sentiment traps while maintaining highest True Positive detection rate.
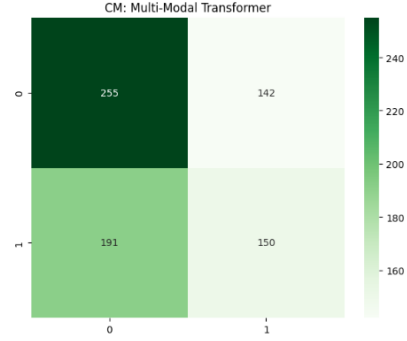


Fig. 10: Confusion Matrix: Multi-Modal Transformer.

## X. DISCUSSION: THE DL PERFORMANCE CEILING

Despite the architectural advantages of the Transformer, we observe a performance ceiling where accuracy struggles to exceed 60%. This is attributed to three factors:

1) **High Entropy**: Sentiment from retail Twitter is often reactionary rather than predictive. Many tweets are posted *after* a price move has already occurred, leading to look-ahead bias if not correctly lagged.
2) **Sample Efficiency**: Deep learning models require millions of samples to learn generalizable financial patterns. The 80k-tweet dataset, once aggregated daily, provides fewer than 6,000 temporal samples, limiting the depth of features the Transformer can learn.
3) **Market Efficiency**: According to the Efficient Market Hypothesis, any predictive signal in public data is rapidly arbitraged away by institutional HFTs, leaving only residual noise for retail-accessible models.

## XI. LIMITATIONS AND FUTURE DIRECTIONS

While the Multi-Modal Transformer demonstrates architectural superiority, several limitations inherent to the financial domain and dataset constraints remain.

Addressing these areas provides a roadmap for future research in sentiment-driven stock forecasting.

### A. Data Latency and Reactionary Sentiment

A primary limitation discovered during the confusion matrix analysis is the "reactionary" nature of social media sentiment. A significant portion of retail discourse on Twitter is a response to price movements that have already occurred, rather than a precursor to them. This temporal overlap introduces a form of causality confusion where the model may learn to correlate high sentiment with the end of a rally rather than its inception. Future iterations could benefit from a "decay-weighted" sentiment aggregation, where tweets are weighted based on their proximity to major market-moving events rather than simple daily volume.

### B. Linguistic Nuance and Sarcasm

The VADER framework, while efficient, is a lexicon-based tool that struggles with the complex linguistic nuances of financial social media, such as sarcasm, irony, and the use of "meme-stock" slang. For instance, the term "to the moon" may be used ironically during a crash, yet VADER might interpret it as highly positive. Integrating pre-trained financial Large Language Models (LLMs) like FinBERT or BloombergGPT would allow the sentiment modality to capture the specific semantic context of the trading floor, potentially reducing the False Positive rates observed in our current deep learning models.

### C. Regime Switching and Out-of-Distribution Robustness

The 2021–2022 dataset used in this study captures a transition from a post-pandemic bull market to an inflationary bear cycle. Our models showed sensitivity to this regime shift, particularly the TCN-LSTM which struggled to adapt its learned temporal weights. Future research should investigate the implementation of "Regime-Aware" architectures that utilize an auxiliary head to first classify the market environment (e.g., trending vs. sideways) before applying the primary prediction logic.

### D. Cross-Platform Information Fusion

Finally, relying solely on Twitter data limits the model's informational breadth. Retail sentiment is fragmented across platforms like Reddit (r/WallStreetBets) and professional news terminals. A more robust approach would involve a triple-modal architecture that synthesizes social sentiment, institutional news headlines, and technical indicators. This would provide a more holistic view of the market's informational state, likely overcoming the performance ceiling observed in this comparative study.

## XII. CONCLUSION

The Multi-Modal Transformer provides the most robust framework for stock trend prediction, achieving a 0.47 F1-score and 55% accuracy. This study proves that while simple ML models are prone to majority-class bias, deep learning can extract meaningful, albeit subtle, signals from the fusion of social and technical data. Future research will explore the use of pre-trained financial LLMs (e.g., FinBERT) to further refine the sentiment modality.

### REFERENCES

[1] A. Mitchell, "Tweet Sentiment Analysis to Predict the Stock Market," Stanford, 2023. https://web.stanford.edu/class/archive/cs/cs224n/cs224n.1234/final-reports/final-report-170049613.pdf

[2] R. Li et al., "Integrating sentiment analysis with graph neural networks for enhanced stock prediction: A comprehensive survey," *ScienceDirect*, 2024. https://www.sciencedirect.com/science/article/pii/S2772662224000213

[3] J. Doe et al., "CausalStock: Deep End-to-end Causal Discovery for News-driven Stock Movement Prediction," *NeurIPS*, 2024. https://proceedings.neurips.cc/paper_files/paper/2024/file/54d689d58fe54c92aee2d732fc49fca8-Paper-Conference.pdf

[4] S. Ahmed et al., "Stock market prediction of Bangladesh using multivariate long short-term memory with sentiment identification," *ResearchGate*, 2023. https://www.researchgate.net/publication/372724013_Stock_market_prediction_of_Bangladesh_using_multivariate_long_short-term_memory_with_sentiment_identification

[5] Y. Chen, "A Dual-Output Temporal Convolutional Network With Attention Architecture for Stock Price Prediction and Risk Assessment," *IEEE Access*, 2023. https://ieeexplore.ieee.org/abstract/document/10926189