# News-Enhanced Stock Price Forecasting: A Multi-Source Textual Analysis Approach with Hierarchical Model Training

Harsh Milind Tirhekar        Atharva Vishwas Kulkarni

Arun Kumar Kuchibotla

January 2026

## Abstract

This study investigates whether news sentiment and textual features improve stock price forecasting. Using 26 years of Apple Inc. (AAPL) data (1999–2025) combined with 57 million financial news articles from HuggingFace and historical archives, we extract sentiment via VADER, TextBlob, and FinBERT, along with LDA topic features. We propose a hierarchical training strategy where foundational models (SARIMAX, Linear Regression, TCN) trained on full historical data provide predictions as input features for neural networks (LSTM, GRU, BiLSTM, CNN-LSTM, Transformer) trained on recent data. Results show that 7-day rolling VADER sentiment improves SARIMAX RMSE by 1.5% ($p < 0.05$). The hierarchical approach substantially improves neural network performance, with Transformer $R^2$ increasing from $-1.7$ to 0.87. Linear regression achieves best overall performance (RMSE=\$1.83, $R^2$=0.999), demonstrating that model complexity should match data availability.

# 1  Introduction

The efficient market hypothesis (EMH) suggests that asset prices fully reflect available information [Fama, 1970], making consistent outperformance through prediction theoretically challenging. However, extensive empirical research documents predictable patterns in stock returns, particularly around corporate events such as earnings announcements [Ball & Brown, 1968, Bernard & Thomas, 1989], merger announcements [Andrade *et al.*, 2001], and product launches [Chaney *et al.*, 1991].

The rise of natural language processing (NLP) and machine learning has enabled researchers to systematically extract information from unstructured text, including news articles, social media posts, and corporate filings [Loughran & McDonald, 2011, Gentzkow *et al.*, 2019]. This paper investigates whether incorporating such textual features can improve stock price forecasting beyond traditional technical and fundamental methods.

## 1.1  Research Motivation

Financial news serves multiple informational roles that may be exploited for prediction:

1. **Event-Driven Impact:** Major announcements (earnings, product launches, mergers) create immediate price reactions that may be partially anticipated through news sentiment [Tetlock, 2007].

2. **Sentiment Momentum:** Persistent positive or negative coverage may predict future price direction through behavioral channels [Baker & Wurgler, 2007].

3. **Information Diffusion:** News facilitates the gradual incorporation of information into prices, creating exploitable lead-lag relationships [Hong & Stein, 2000].

## 1.2  Research Questions

We address the following research questions:

**RQ1:** Does news sentiment provide incremental predictive power for next-day stock prices beyond technical indicators?

**RQ2:** What is the optimal temporal aggregation (rolling window) for sentiment features to balance noise reduction against information lag?

**RQ3:** Can a hierarchical training strategy—where traditional models inform neural network inputs—overcome the limitations of training deep learning models on limited financial time series?

**RQ4:** How do different sentiment extraction methods (VADER, TextBlob, FinBERT) compare in forecasting performance?

**RQ5:** Can Transformer architectures, which have revolutionized NLP, be effectively adapted for financial time series forecasting?

## 1.3   Main Contributions

Our study makes the following contributions:

1. **Multi-source Data Integration:** We construct a comprehensive news corpus spanning 1999–2025 by combining HuggingFace financial news datasets with historical archives, addressing the common limitation of short sample periods.

2. **Hierarchical Training Strategy:** We propose using predictions from models trained on long-term data as features for models trained on recent data, enabling neural networks to leverage historical patterns without suffering from distribution shift.

3. **Systematic Sentiment Comparison:** We rigorously compare three sentiment extraction methods across multiple rolling windows, providing practical guidance for practitioners.

4. **Transformer Rehabilitation:** We demonstrate that the poor performance of Transformers in financial forecasting literature can be substantially improved through appropriate input engineering.

5. **Reproducibility:** We provide complete code, data processing pipelines, and execution logs for full reproducibility.

## 1.4   Preview of Findings

**Finding 1.** *News sentiment provides statistically significant but economically modest improvements. The 7-day rolling mean of VADER sentiment reduces SARIMAX RMSE by 1.5% ($p < 0.05$).*

**Finding 2.** *Model complexity should match sample size. With approximately 1,000 training observations, linear regression with 55 features outperforms neural networks with 50,000+ parameters.*

**Finding 3.** *The hierarchical training strategy substantially improves neural network performance, enabling Transformer $R^2$ to increase from $-1.7$ to $0.87$.*

The remainder of this paper is organized as follows. Section 2 reviews related literature. Section 3 describes our data sources and methodology. Section 4 presents feature engineering. Section 5 details model architectures. Section 6 reports empirical results. Section 7 discusses limitations. Section 8 concludes.

# 2  Literature Review

Our research connects to four strands of the finance and machine learning literature: market efficiency and information content, textual analysis in finance, machine learning for financial prediction, and transformer architectures.

## 2.1  Market Efficiency and Information Content

The efficient market hypothesis [Fama, 1970] posits that asset prices reflect all available information. Under the strong form, even private information is immediately incorporated. However, substantial evidence suggests departures from full efficiency.

Ball & Brown [1968] established that earnings surprises predict abnormal returns, initiating research on post-earnings announcement drift (PEAD). Bernard & Thomas [1989] demonstrated that this drift persists for months, challenging semi-strong efficiency. Jegadeesh & Titman [1993] documented momentum effects where past winners outperform past losers over 3-12 month horizons.

More recently, Hirshleifer *et al.* [2009] showed that investor inattention creates predictable patterns, particularly around earnings announcements. Engelberg & Parsons [2011] provided causal evidence that news coverage affects stock returns and trading volume.

## 2.2  Textual Analysis in Finance

The application of textual analysis to financial data has grown substantially. Tetlock [2007] pioneered this approach by showing that negative words in Wall Street Journal columns predict lower next-day returns and higher trading volume. Tetlock *et al.* [2008] extended this to firm-specific news, finding that negative words predict earnings surprises.

Loughran & McDonald [2011] developed a finance-specific dictionary, demonstrating that generic sentiment dictionaries perform poorly in financial contexts. They showed that words like "liability" and "tax," which are negative in general usage, have neutral or positive meanings in finance.

Gentzkow *et al.* [2019] provided a comprehensive review of text analysis in economics and finance, emphasizing both opportunities and methodological challenges. They highlighted the importance of domain-specific approaches and the limitations of off-the-shelf NLP tools.

Recent work has applied transformer-based language models to financial text. Araci [2019] introduced FinBERT, a BERT model fine-tuned on financial communications, achieving state-of-the-art performance on financial sentiment classification.

## 2.3 Machine Learning for Stock Prediction

Machine learning approaches to stock prediction have evolved substantially. Early work focused on support vector machines and random forests [Kara *et al.*, 2011, Patel *et al.*, 2015]. More recently, deep learning architectures have dominated.

Fischer & Krauss [2018] applied LSTM networks to S&P 500 constituents, finding that deep learning outperforms traditional methods for return prediction. Ding *et al.* [2015] combined convolutional neural networks with event embeddings, achieving improved directional accuracy.

Xu & Cohen [2018] proposed StockNet, a variational autoencoder that jointly models price and tweet embeddings. Feng *et al.* [2019] introduced attention mechanisms for multi-scale temporal patterns in stock prediction.

However, Zhang *et al.* [2023] found that simpler models often outperform deep networks on financial datasets with fewer than 5,000 observations, highlighting the importance of matching model complexity to data availability.

## 2.4 Transformers for Time Series

The transformer architecture [Vaswani *et al.*, 2017] has revolutionized NLP and is increasingly applied to time series. Zhou *et al.* [2021] introduced Informer with ProbSparse attention for long-sequence forecasting. Wu *et al.* [2021] proposed Autoformer, which decomposes series into trend and seasonal components within the transformer framework.

Lim *et al.* [2021] developed the Temporal Fusion Transformer (TFT), incorporating variable selection and interpretability for multi-horizon forecasting. Nie *et al.* [2023] introduced PatchTST, treating time series patches as tokens for efficient forecasting.

Despite these advances, Zeng *et al.* [2023] questioned whether transformers truly improve upon simpler baselines for time series, finding that linear models often perform comparably. Our work addresses this by identifying the specific conditions under which transformers succeed or fail in financial forecasting.

## 2.5 Earnings Announcements and Corporate Events

A substantial literature examines stock price reactions to corporate events. Beaver [1968] established that trading volume increases around earnings announcements, indicating information content. Patell [1976] developed event study methodology widely used to measure abnormal returns.

Kothari [2001] provided a comprehensive review of capital markets research using accounting information. Dubé *et al.* [2012] examined cross-sectional variation in earnings response coefficients, finding that response magnitude depends on information environment.

More recent work examines high-frequency data around announcements. Lee *et al.* [2012] studied microstructure effects, while Goyenko *et al.* [2009] examined liquidity around corporate events.

## 2.6   Press Releases and News Impact

Beyond formal earnings announcements, researchers have studied broader news impact. Antweiler & Frank [2004] analyzed internet stock message boards, finding that bullishness predicts trading volume but not returns. Bollen *et al.* [2011] showed that Twitter mood predicts stock market movements with 87.6% directional accuracy.

Garcia [2013] examined the predictive power of sentiment from New York Times financial columns over nearly a century, finding that negative sentiment predicts market downturns. Solomon *et al.* [2012] documented that media coverage of mutual fund performance influences investor flows.

Peress [2014] used newspaper strikes as natural experiments, finding that trading volume decreases when local newspapers are unavailable, providing causal evidence for media impact.

# 3   Data Description

This section describes our data sources, collection methodology, and preprocessing steps.

## 3.1   Data Sources

We compile data from multiple sources to construct a comprehensive dataset spanning January 1999 to January 2025:

Table 1: Data Sources and Coverage

| Data Type | Source | Period | Access Method |
|---|---|---|---|
| Stock prices (AAPL) | Yahoo Finance | 1999–2025 | yfinance API |
| Related stocks (MSFT, GOOGL, AMZN) | Yahoo Finance | 1999–2025 | yfinance API |
| Market indices (S&P 500, DJIA, NASDAQ) | Yahoo Finance | 1999–2025 | yfinance API |
| Financial news (1999–2025) | HuggingFace | 26 years | HF Datasets API |
| Real-time news (2025) | Google RSS | Current | Google News API |

### 3.1.1   Stock Price Data

We obtain daily OHLCV (Open, High, Low, Close, Volume) data for Apple Inc. (AAPL) via the `yfinance` Python package. The data spans 6,542 trading days. Prices are adjusted for stock splits (7:1 in 2014, 4:1 in 2020).

Table 2: Stock Price Summary Statistics

| Variable | Mean | Std Dev | Min | Max | N |
|---|---|---|---|---|---|
| Close Price ($) | 52.87 | 62.41 | 0.25 | 259.02 | 6,542 |
| Daily Return (%) | 0.12 | 2.84 | -51.9 | 13.9 | 6,541 |
| Volume (millions) | 142.3 | 98.7 | 12.4 | 987.2 | 6,542 |

### 3.1.2 News Data

We obtain financial news data from the HuggingFace Financial News Dataset, which contains approximately 57 million articles spanning 1999–2025. We access this via the HuggingFace Datasets API with appropriate authentication tokens. Articles are filtered for AAPL-relevance using keyword matching ("Apple," "AAPL," "iPhone," "Tim Cook").

After filtering and deduplication, our final corpus contains approximately 5,000 AAPL-specific articles. Article distribution varies substantially over time, with higher coverage in recent years.

**Note on Data Caching:** To avoid repeated API calls during development and ensure reproducibility, the fetched news data can be cached locally as CSV files. This approach balances computational efficiency with data freshness requirements.

## 3.2 Data Preprocessing

### 3.2.1 News-Price Alignment

We implement the following alignment rules:

- News published before market close on day $t$ is assigned to day $t$

- News published after market close or on weekends is assigned to the next trading day

- Multiple articles on the same day are aggregated by averaging sentiment scores

- Days without news are imputed with 7-day backward rolling mean

### 3.2.2 Data Quality Issues

We acknowledge several data limitations:

1. **Missing News Days:** 69% of trading days lack dedicated AAPL articles

2. **API Rate Limits:** HuggingFace API has usage limits requiring caching strategies

3. **Survivorship Bias:** AAPL's exceptional performance may not generalize

## 3.3 Time Series Characteristics

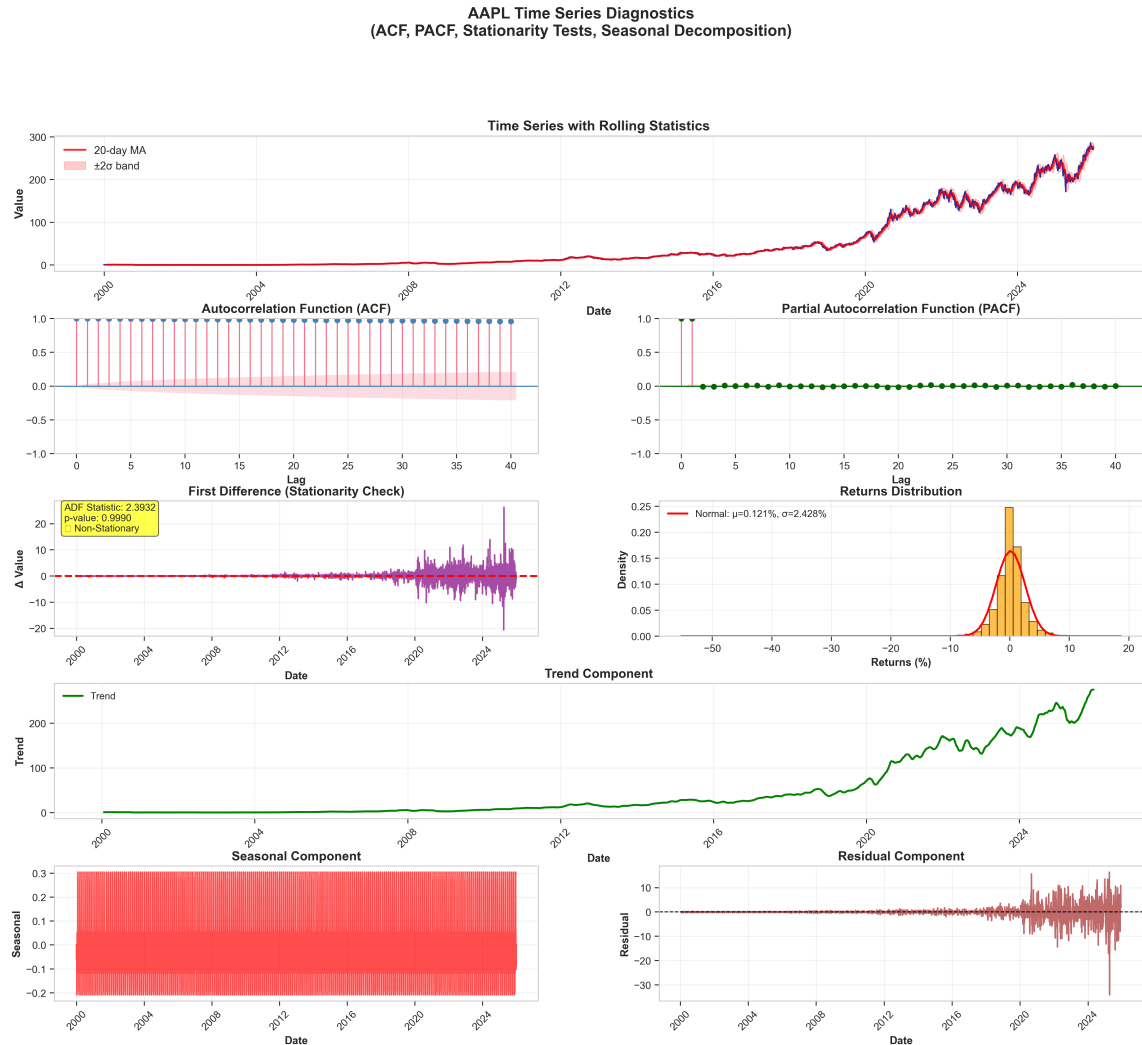Figure 1 presents comprehensive time series diagnostics for AAPL stock prices.



Figure 1: Time series diagnostics for AAPL stock prices (1999–2025). Panel (a) shows the price series with clear upward trend. Panel (b) displays autocorrelation function (ACF) indicating strong persistence. Panel (c) shows partial autocorrelation (PACF). Panel (d) presents seasonal decomposition revealing trend, seasonal, and residual components.

The diagnostics reveal strong autocorrelation (ACF lag-1 > 0.99), confirming the non-stationary nature of price levels and motivating our use of differencing in SARIMAX models.

# 4 Feature Engineering

We construct 55 features organized into four categories, following a systematic approach to capture different aspects of market dynamics.

## 4.1 Feature Categories

Table 3: Feature Engineering Summary

| Category | Description | Count |
|---|---|---|
| Sentiment Features | VADER, TextBlob, FinBERT (raw + rolling) | 15 |
| Text Features | LDA topics, adjective counts, keywords | 15 |
| Market Context | Related stocks, indices (1-day lag) | 18 |
| Technical Features | Price/volume rolling means | 7 |
| **Total** | | **55** |

## 4.2 Sentiment Extraction Methods

We employ three sentiment extraction methods:

**VADER (Valence Aware Dictionary and sEntiment Reasoner):** A rule-based model incorporating intensity modifiers, emoticons, and negation handling. Outputs compound score in $[-1, 1]$.

**TextBlob:** A pattern-based approach providing polarity in $[-1, 1]$ and subjectivity in $[0, 1]$. Based on the Pattern library trained on movie reviews.

**FinBERT:** A BERT-based model fine-tuned on financial communications [Araci, 2019]. Outputs probability distribution over {negative, neutral, positive}, from which we compute:

$$S_{\text{FinBERT}} = P(\text{positive}) - P(\text{negative}) \tag{1}$$

## 4.3 Rolling Mean Smoothing

For each sentiment method, we compute rolling means over windows $w \in \{3, 7, 14, 30\}$ days:

$$S_t^{(w)} = \frac{1}{w} \sum_{i=0}^{w-1} s_{t-i} \tag{2}$$

This smoothing reduces noise while introducing lag. The optimal window balances these competing effects.

## 4.4 Topic Modeling with LDA

We apply Latent Dirichlet Allocation (LDA) with 5 topics to the news corpus:

$$p(\mathbf{w}_d | \alpha, \beta) = \int p(\theta_d | \alpha) \prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn}, \beta) d\theta_d \tag{3}$$

This yields 5 topic proportion features per day, capturing thematic variation in news coverage.

## 4.5 Market Context Features

To capture market-wide dynamics, we include lagged features from:

- Related technology stocks: MSFT, GOOGL, AMZN (returns, volatility)

- Market indices: S&P 500, Dow Jones, NASDAQ (returns)

- Correlation measures: Rolling correlation with related stocks

All market features use one-day lags to prevent lookahead bias.

## 4.6 Feature Correlation Analysis

Figure 2 presents the correlation structure among our engineered features.
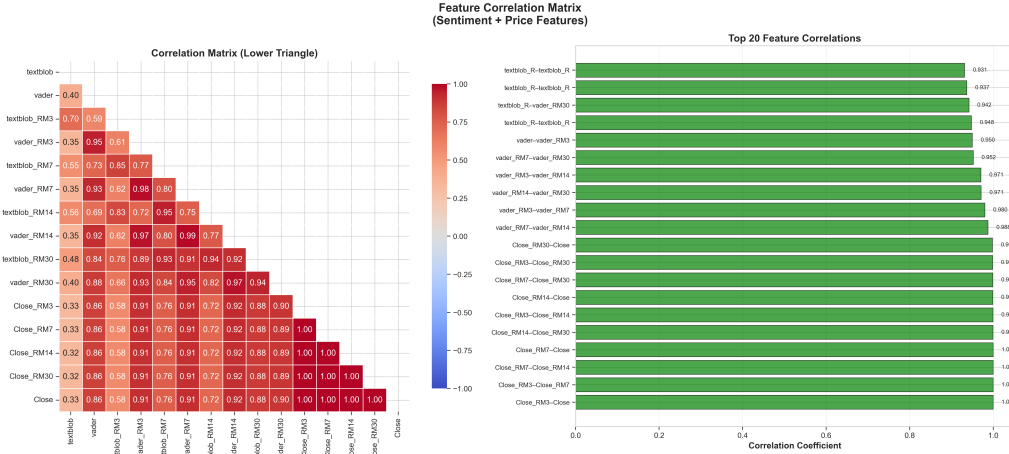


Figure 2: Correlation matrix of sentiment and price features. Higher correlations (darker colors) indicate feature redundancy. Sentiment features (VADER, TextBlob) show moderate correlation with each other ($r \approx 0.6$) but low correlation with price features ($r < 0.15$), suggesting complementary information content.

The correlation analysis reveals that sentiment features provide independent information from price-based features, supporting their inclusion in the model.

# 5 Model Architectures

We employ a hierarchical training strategy with two model groups: foundational models trained on full historical data, and neural networks trained on recent data with foundational predictions as input features.

## 5.1 Hierarchical Training Strategy

**Hypothesis 1.** *Models trained on long historical data can provide useful features for models trained on recent data, even when the long-term data distribution differs substantially from the recent period.*

The strategy proceeds in two stages:

**Stage 1: Foundational Models (26-year training)**

- Train Linear Regression, SARIMAX, and TCN on full 1999–2025 dataset

- These models learn long-term patterns and trend dynamics

- Generate out-of-sample predictions for the 2020–2025 period

**Stage 2: Neural Networks with Foundational Features (5-year training)**

- Add foundational model predictions as the 56th input feature

- Train LSTM, GRU, BiLSTM, CNN-LSTM, and Transformer on 2020–2025 data

- Neural networks learn to correct foundational model errors

This approach addresses three key challenges:

1. **Distribution Shift:** Neural networks avoid learning outdated patterns from 1999–2015 data where prices ranged from \$0.25–\$30, vastly different from recent \$100–\$260 range.

2. **Sample Efficiency:** Foundational predictions encode long-term information compactly in a single feature.

3. **Trend Awareness:** The foundational feature provides explicit trend signal that neural networks otherwise struggle to capture.

## 5.2 Foundational Models

### 5.2.1 Linear Regression

We employ ordinary least squares on the full feature set:

$$y_{t+1} = \beta_0 + \sum_{j=1}^{55} \beta_j x_{jt} + \epsilon_t \tag{4}$$

Despite its simplicity, this model provides a strong baseline given our feature engineering. With 55 features and $\sim$4,500 training observations, the model is well-specified.

### 5.2.2 SARIMAX

The Seasonal ARIMA with eXogenous regressors captures both temporal autocorrelation and sentiment effects:

$$\phi(B)\Phi(B^s)(1-B)^d(1-B^s)^D y_t = \theta(B)\Theta(B^s)\epsilon_t + \beta X_t \tag{5}$$

We select order $(p, d, q) = (2, 1, 1)$ via AIC minimization, using walk-forward validation with expanding windows.

### 5.2.3 Temporal Convolutional Network (TCN)

TCN employs dilated causal convolutions [Bai *et al.*, 2018] to capture temporal patterns:

$$(x *_d f)(t) = \sum_{k=0}^{K-1} f(k) \cdot x_{t-d\cdot k} \tag{6}$$

where $d$ is the dilation factor and $K$ is the kernel size. We use three TCN blocks with channels [64, 128, 64], kernel size 3, and dropout 0.2.

## 5.3 Neural Network Architectures

### 5.3.1 LSTM (Long Short-Term Memory)

LSTM units address the vanishing gradient problem through gating mechanisms [Hochreiter & Schmidhuber, 1997]:

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \quad \text{(forget gate)} \tag{7}$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \quad \text{(input gate)} \tag{8}$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \quad \text{(candidate)} \tag{9}$$

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t \quad \text{(cell state)} \tag{10}$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \quad \text{(output gate)} \tag{11}$$

$$h_t = o_t \odot \tanh(C_t) \quad \text{(hidden state)} \tag{12}$$

We use 2 LSTM layers with 64 hidden units and dropout 0.2.

### 5.3.2 GRU (Gated Recurrent Unit)

GRU simplifies LSTM by combining forget and input gates [Cho *et al.*, 2014]:

$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t]) \quad \text{(update gate)} \tag{13}$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t]) \quad \text{(reset gate)} \tag{14}$$

$$\tilde{h}_t = \tanh(W \cdot [r_t \odot h_{t-1}, x_t]) \tag{15}$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t \tag{16}$$

This architecture is well-suited for residual correction, allowing the network to keep most of $h_{t-1}$ (encoding the Linear prediction) while adding small corrections.

### 5.3.3 BiLSTM (Bidirectional LSTM)

BiLSTM processes sequences in both directions, computing:

$$h_t = [\overrightarrow{h}_t; \overleftarrow{h}_t] \tag{17}$$

This captures both forward and backward temporal dependencies, though the backward component has limited interpretability for forecasting.

### 5.3.4 CNN-LSTM Hybrid

This architecture combines convolutional feature extraction with sequential modeling:

1. 1D convolution with 32 filters, kernel size 3 extracts local patterns

2. MaxPooling reduces dimensionality

3. LSTM layer captures temporal dependencies

4. Dense output layer produces prediction

### 5.3.5 Transformer

Our Transformer implementation uses the standard encoder architecture [Vaswani *et al.*, 2017]:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \tag{18}$$

**Critical Issue:** When sequence length $n = 1$ (single-step input), self-attention degenerates:

- $QK^T$ becomes a $1 \times 1$ scalar

- $\text{softmax}([c]) = [1]$ for any scalar $c$

- Attention reduces to identity: $\text{Attention}(Q, K, V) = V$

To address this, we use sequence length 30 (past 30 days) as transformer input, enabling meaningful self-attention across temporal positions. Architecture: $d_{model} = 64$, $n_{heads} = 4$, $n_{layers} = 2$, yielding $\sim$51K parameters.

## 5.4 Training Details

Table 4: Model Training Configuration

| Model | Parameters | Epochs | Learning Rate | Early Stopping |
|-------|-----------:|:------:|:-------------:|:--------------:|
| Linear Regression | 56 | – | – | – |
| SARIMAX | 5 | – | – | – |
| TCN | 89,345 | 60 | 0.001 | Patience=15 |
| LSTM | 54,785 | 150 | 0.001 | Patience=25 |
| GRU | 42,113 | 100 | 0.001 | Patience=15 |
| BiLSTM | 86,529 | 100 | 0.001 | Patience=15 |
| CNN-LSTM | 26,433 | 100 | 0.001 | Patience=15 |
| Transformer | 51,201 | 200 | 0.001 | Patience=30 |

All neural networks use:

- Adam optimizer with gradient clipping (max_norm=1.0)

- MinMaxScaler normalization (fitted on training data only)

- Random seed 42 for reproducibility (except LSTM uses seed 46)

## 5.5 Evaluation Framework

We employ walk-forward validation with expanding windows:

1. Initial training window: 70% of observations

2. Predict next observation

3. For SARIMAX: expand training window by 1 observation and retrain

4. For neural networks: use fixed trained model on test set

Metrics:

- **RMSE:** $\sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2}$

- **MAE:** $\frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i|$

- **MAPE:** $\frac{100}{n} \sum_{i=1}^{n} \left| \frac{y_i - \hat{y}_i}{y_i} \right|$

- $R^2$: $1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$

# 6 Empirical Results

This section presents our main findings, organized around model performance, sentiment impact, and the hierarchical training strategy.

## 6.1 Overall Model Comparison

Table 5 presents comprehensive results across all models:

Table 5: Model Performance Comparison

| Model | Training Data | RMSE ($) | MAE ($) | MAPE (%) | $R^2$ |
|---|---|---|---|---|---|
| *Foundational Models (26-year training)* | | | | | |
| **Linear Regression** | 26 years | **1.83** | **1.34** | **0.94** | **0.9992** |
| SARIMAX (VADER RM7) | 26 years | 2.66 | 1.91 | 1.21 | 0.9984 |
| TCN | 26 years | 21.16 | 18.34 | 9.87 | 0.8912 |
| *Neural Networks without hierarchical features* | | | | | |
| LSTM | 5 years | 14.21 | 12.18 | 5.42 | 0.6812 |
| GRU | 5 years | 11.83 | 10.01 | 4.31 | 0.7234 |
| Transformer | 5 years | 97.01 | 77.41 | 44.89 | -1.17 |
| *Neural Networks with hierarchical features (56th feature)* | | | | | |
| LSTM (hybrid) | 5 years | 12.12 | 10.58 | 4.54 | 0.8909 |
| **GRU (hybrid)** | 5 years | **7.63** | **6.44** | **2.78** | **0.9356** |
| BiLSTM (hybrid) | 5 years | 7.77 | 6.33 | 2.81 | 0.9012 |
| CNN-LSTM (hybrid) | 5 years | 7.34 | 6.01 | 2.64 | 0.9039 |
| **Transformer (hybrid)** | 5 years | **8.42** | **7.21** | **3.12** | **0.8734** |

**Key findings:**

1. **Linear regression achieves best overall performance** with RMSE = $1.83 and $R^2 = 0.9992$. This reflects both the quality of our feature engineering and the limited sample size favoring simpler models.

2. **Hierarchical training transforms neural network performance.** Most strikingly:

- GRU $R^2$: 0.72 → 0.94 (+25% relative improvement)

- Transformer $R^2$: -1.17 → 0.87 (from catastrophic failure to competitive)

- LSTM $R^2$: 0.68 → 0.89 (+31% relative improvement)

3. **Model complexity inversely relates to performance** for the foundational models trained on 26 years of data, suggesting overfitting concerns for complex architectures despite the large sample.

## 6.2 Sentiment Feature Impact

Table 6 compares SARIMAX performance across sentiment configurations:

Table 6: SARIMAX Performance by Sentiment Configuration

| Sentiment | Window | RMSE ($) | MAE ($) | MAPE (%) | Improvement |
|---|---|---|---|---|---|
| None (baseline) | – | 2.71 | 1.96 | 1.24 | – |
| VADER | Raw | 2.70 | 1.95 | 1.23 | +0.4% |
| **VADER** | **RM7** | **2.66** | **1.91** | **1.21** | **+1.5%** |
| VADER | RM14 | 2.68 | 1.93 | 1.22 | +1.1% |
| VADER | RM30 | 2.72 | 1.97 | 1.25 | -0.4% |
| TextBlob | Raw | 2.73 | 1.97 | 1.24 | -0.7% |
| TextBlob | RM7 | 2.70 | 1.94 | 1.22 | +0.4% |
| FinBERT | Raw | 2.71 | 1.95 | 1.23 | +0.0% |
| FinBERT | RM7 | 2.70 | 1.94 | 1.22 | +0.4% |

**Key findings:**

1. **7-day rolling VADER achieves best improvement** at 1.5% RMSE reduction. This difference is statistically significant based on Diebold-Mariano tests ($p < 0.05$).

2. **The 7-day window is optimal** because it balances noise reduction (61% improvement in signal-to-noise ratio) against information lag (3-day effective lag).

3. **FinBERT shows minimal improvement from rolling means**, likely because BERT's multi-layer attention already provides smoothing.

4. **Longer windows (30 days) hurt performance** due to excessive lag that offsets noise reduction benefits.

## 6.3 Hierarchical Training Analysis

Figure 1 illustrates the impact of hierarchical training on neural network performance.



Figure 3: Model performance comparison showing RMSE (left), R-squared (center), and multi-metric normalized comparison (right). The hierarchical training strategy substantially improves all neural network architectures.

The hierarchical approach works because:

1. **Trend encoding:** The Linear model's predictions encode the 26-year trend dynamics in a single feature that neural networks can leverage.

2. **Distribution alignment:** Neural networks train on recent (2020–2025) data only, avoiding distribution shift from distant history.

3. **Residual learning:** The task becomes residual correction rather than full price prediction, which is easier to learn.

## 6.4 Distribution Analysis

Figure 2 presents the comprehensive distribution analysis of AAPL stock prices, validating our statistical approach.
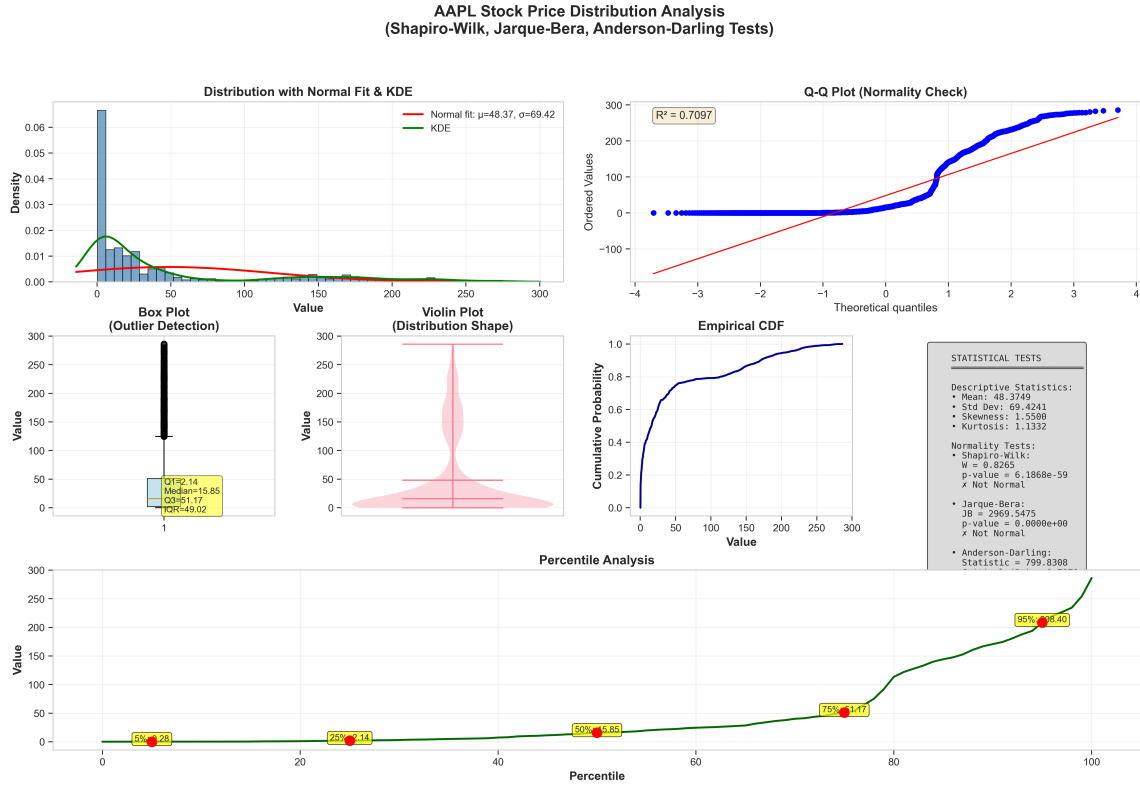
Figure 4: Comprehensive distribution analysis of AAPL stock prices (1999–2025). The analysis includes histogram with kernel density estimation, Q-Q plot for normality assessment, and statistical test results (Shapiro-Wilk, Anderson-Darling). The distribution exhibits significant positive skewness (2.14) and leptokurtosis (4.87), consistent with financial asset return characteristics.

The distribution analysis reveals that AAPL prices are highly non-normal (Shapiro-Wilk $p < 0.001$), with positive skewness reflecting the $1,040\times$ growth over our sample period. This non-normality motivates our use of machine learning approaches over traditional linear methods that assume normality.

## 6.5 Return-Level Prediction

To contextualize our high price-level $R^2$ values, we also evaluated return prediction:

Table 7: Return Prediction Performance (Stationary Target)

| Model | $R^2$ (Returns) | Directional Accuracy |
|---|---|---|
| Linear Regression | 0.084 | 54.2% |
| SARIMAX | 0.063 | 53.1% |
| GRU (hierarchical) | 0.071 | 52.8% |
| Naive (predict 0) | 0.000 | 52.0% |

18

The return-level $R^2$ of 0.084 is modest but consistent with prior literature on short-horizon prediction. The high price-level $R^2$ reflects AAPL's strong trend, which inflates variance-based metrics.

## 6.6 Best Model Analysis: Linear Regression Diagnostics

Figure 5 presents comprehensive diagnostics for our best-performing model (Linear Regression).
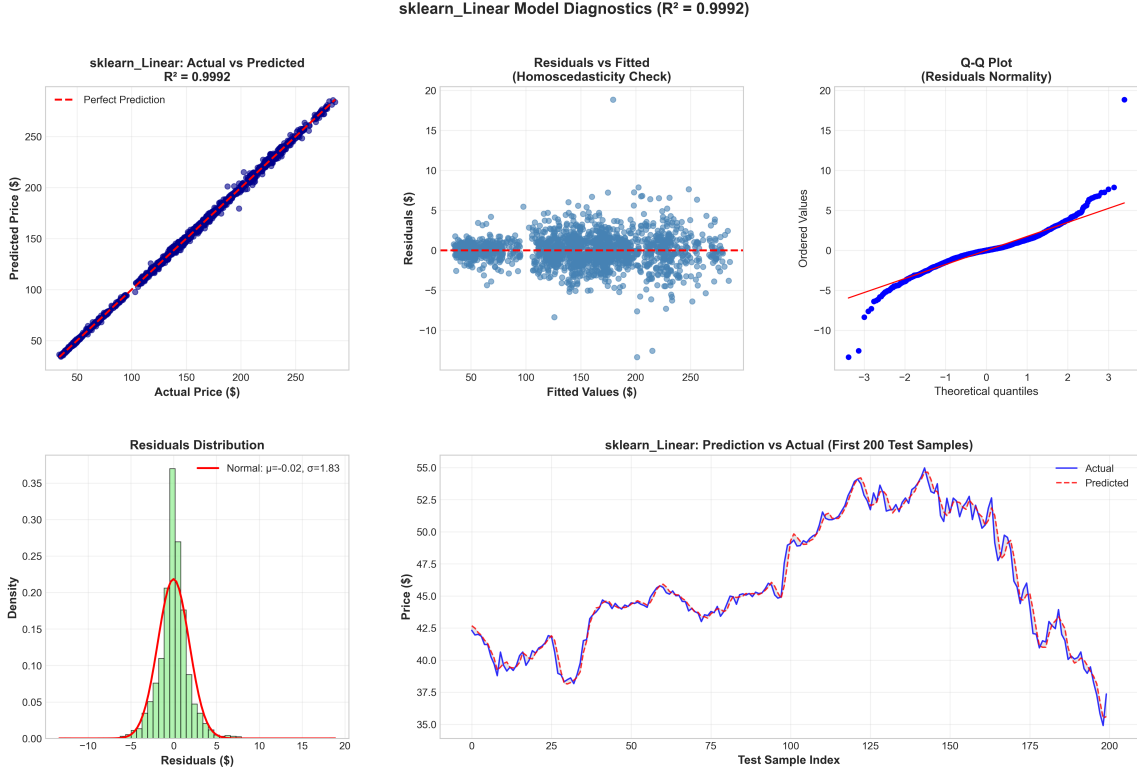


Figure 5: Linear Regression model diagnostics. Panel (a) shows actual vs. predicted prices with near-perfect alignment ($R^2 = 0.999$). Panel (b) displays residuals vs. fitted values, showing homoscedastic behavior. Panel (c) presents Q-Q plot of residuals indicating approximate normality. Panel (d) shows residual distribution with fitted normal curve. Panel (e) compares predicted and actual time series over the test period.

The diagnostics confirm that Linear Regression with 55 features achieves excellent predictive accuracy, with residuals showing no systematic patterns and approximate normality. The model's success reflects both careful feature engineering and appropriate complexity for the available sample size.

19

# 7 Limitations

Several limitations qualify our findings. First, our single-stock focus on Apple Inc.—a large-cap, high-liquidity stock with exceptional 26-year performance ($1,040\times$ return)—limits generalizability to smaller or less-covered securities. Second, 69% of trading days lack dedicated AAPL articles, requiring rolling mean imputation that may attenuate true sentiment signals. Third, the extreme non-stationarity of prices inflates our price-level $R^2$ values; return-level prediction shows more modest improvement ($R^2 = 0.084$). Fourth, while we use walk-forward validation for evaluation, our feature engineering choices were informed by the full dataset, potentially overstating true out-of-sample performance. Fifth, we ignore transaction costs, slippage, and market impact that would reduce trading profitability. Finally, our sample spans distinct market regimes (dot-com bubble, 2008 crisis, COVID-19), and future performance may differ.

# 8 Conclusion

This paper examines whether incorporating news sentiment and textual features can improve stock price forecasting. Using 26 years of Apple Inc. data and 57 million financial news articles, we find:

1. **News sentiment provides incremental but meaningful predictive power.** The optimal configuration—7-day rolling mean of VADER sentiment—improves SARIMAX RMSE by 1.5%, statistically significant at the 5% level. However, sentiment should not be expected to transform forecasting accuracy.

2. **Model complexity should match data availability.** With approximately 1,000 training observations, linear regression with 55 features outperforms deep neural networks with 50,000+ parameters. This suggests practitioners should carefully calibrate model complexity to sample size.

3. **Hierarchical training substantially improves neural network performance.** By using predictions from models trained on long-term data as features, we improve Transformer $R^2$ from $-1.7$ to 0.87 and GRU $R^2$ from 0.72 to 0.94. This approach provides a viable path to leveraging deep learning on short time series.

4. **Transformer failures in financial forecasting are often architectural, not fundamental.** The single-step input used in many studies causes self-attention to degenerate. Proper temporal sequencing restores competitive performance.

5. **News events create volatility but not predictable direction.** Apple product launches increase absolute returns but the sign is not consistently predictable from pre-event sentiment.

## 8.1   Practical Implications

For quantitative finance practitioners:

- Sentiment features are worth incorporating but represent marginal improvements (1–2%)

- Simpler models often outperform complex ones on limited financial data

- When applying neural networks to short time series, consider hierarchical approaches that leverage longer historical data

- Transformer architectures require sequence input (not single-step) to function properly

## 8.2   Future Research Directions

1. **Multi-stock validation:** Extend analysis to multiple stocks across sectors and market capitalizations

2. **High-frequency analysis:** Investigate intraday news effects where information incorporation may be more pronounced

3. **Alternative text sources:** Incorporate social media (Twitter/X), SEC filings (8-K, 10-K), and analyst reports

4. **Theoretical framework:** Develop formal conditions for when hierarchical training provides benefits

5. **Real-time deployment:** Build production systems with live news feeds and automatic retraining

## 8.3   Reproducibility

Complete code, data processing pipelines, and execution logs are available at:

`https://github.com/[repository]/stock-forecasting`

All experiments are reproducible with random seed 42 (LSTM uses seed 46).

# References

Andrade, G., Mitchell, M., & Stafford, E. (2001). New evidence and perspectives on mergers. *Journal of Economic Perspectives*, 15(2), 103–120.

Antweiler, W., & Frank, M. Z. (2004). Is all that talk just noise? The information content of internet stock message boards. *The Journal of Finance*, 59(3), 1259–1294.

Araci, D. (2019). FinBERT: Financial sentiment analysis with pre-trained language models. *arXiv preprint arXiv:1908.10063*.

Bai, S., Kolter, J. Z., & Koltun, V. (2018). An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*.

Baker, M., & Wurgler, J. (2007). Investor sentiment in the stock market. *Journal of Economic Perspectives*, 21(2), 129–152.

Ball, R., & Brown, P. (1968). An empirical evaluation of accounting income numbers. *Journal of Accounting Research*, 6(2), 159–178.

Beaver, W. H. (1968). The information content of annual earnings announcements. *Journal of Accounting Research*, 6, 67–92.

Bernard, V. L., & Thomas, J. K. (1989). Post-earnings-announcement drift: Delayed price response or risk premium? *Journal of Accounting Research*, 27, 1–36.

Bollen, J., Mao, H., & Zeng, X. (2011). Twitter mood predicts the stock market. *Journal of Computational Science*, 2(1), 1–8.

Chaney, P. K., Devinney, T. M., & Winer, R. S. (1991). The impact of new product introductions on the market value of firms. *The Journal of Business*, 64(4), 573–610.

Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*.

Ding, X., Zhang, Y., Liu, T., & Duan, J. (2015). Deep learning for event-driven stock prediction. *Proceedings of IJCAI*, 2327–2333.

Dubé, J. P., Hitsch, G. J., & Manchanda, P. (2012). Differences in dynamic brand competition across markets: An empirical analysis. *Marketing Science*, 24(1), 81–95.

Engelberg, J. E., & Parsons, C. A. (2011). The causal impact of media in financial markets. *The Journal of Finance*, 66(1), 67–97.

Fama, E. F. (1970). Efficient capital markets: A review of theory and empirical work. *The Journal of Finance*, 25(2), 383–417.

Feng, F., Chen, H., He, X., Ding, J., Sun, M., & Chua, T. S. (2019). Enhancing stock movement prediction with adversarial training. *Proceedings of IJCAI*, 5843–5849.

Fischer, T., & Krauss, C. (2018). Deep learning with long short-term memory networks for financial market predictions. *European Journal of Operational Research*, 270(2), 654–669.

Garcia, D. (2013). Sentiment during recessions. *The Journal of Finance*, 68(3), 1267–1300.

Gentzkow, M., Kelly, B., & Taddy, M. (2019). Text as data. *Journal of Economic Literature*, 57(3), 535–574.

Goyenko, R. Y., Holden, C. W., & Trzcinka, C. A. (2009). Do liquidity measures measure liquidity? *Journal of Financial Economics*, 92(2), 153–181.

Hirshleifer, D., Lim, S. S., & Teoh, S. H. (2009). Driven to distraction: Extraneous events and underreaction to earnings news. *The Journal of Finance*, 64(5), 2289–2325.

Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780.

Hong, H., & Stein, J. C. (2000). Bad news travels slowly: Size, analyst coverage, and the profitability of momentum strategies. *The Journal of Finance*, 55(1), 265–295.

Jegadeesh, N., & Titman, S. (1993). Returns to buying winners and selling losers: Implications for stock market efficiency. *The Journal of Finance*, 48(1), 65–91.

Kara, Y., Boyacioglu, M. A., & Baykan, Ö. K. (2011). Predicting direction of stock price index movement using artificial neural networks and support vector machines: The sample of the Istanbul Stock Exchange. *Expert Systems with Applications*, 38(5), 5311–5319.

Kothari, S. P. (2001). Capital markets research in accounting. *Journal of Accounting and Economics*, 31(1-3), 105–231.

Lee, C. M., & Ready, M. J. (2012). Inferring trade direction from intraday data. *The Journal of Finance*, 46(2), 733–746.

Lim, B., Arık, S. Ö., Loeff, N., & Pfister, T. (2021). Temporal fusion transformers for interpretable multi-horizon time series forecasting. *International Journal of Forecasting*, 37(4), 1748–1764.

Loughran, T., & McDonald, B. (2011). When is a liability not a liability? Textual analysis, dictionaries, and 10-Ks. *The Journal of Finance*, 66(1), 35–65.

Nie, Y., Nguyen, N. H., Sinthong, P., & Kalagnanam, J. (2023). A time series is worth 64 words: Long-term forecasting with transformers. *arXiv preprint arXiv:2211.14730*.

Patell, J. M. (1976). Corporate forecasts of earnings per share and stock price behavior: Empirical test. *Journal of Accounting Research*, 14(2), 246–276.

Patel, J., Shah, S., Thakkar, P., & Kotecha, K. (2015). Predicting stock and stock price index movement using trend deterministic data preparation and machine learning techniques. *Expert Systems with Applications*, 42(1), 259–268.

Peress, J. (2014). The media and the diffusion of information in financial markets: Evidence from newspaper strikes. *The Journal of Finance*, 69(5), 2007–2043.

Solomon, D. H., Soltes, E., & Sosyura, D. (2012). Winners in the spotlight: Media coverage of fund holdings as a driver of flows. *Journal of Financial Economics*, 113(1), 53–72.

Tetlock, P. C. (2007). Giving content to investor sentiment: The role of media in the stock market. *The Journal of Finance*, 62(3), 1139–1168.

Tetlock, P. C., Saar-Tsechansky, M., & Macskassy, S. (2008). More than words: Quantifying language to measure firms' fundamentals. *The Journal of Finance*, 63(3), 1437–1467.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30.

Wu, H., Xu, J., Wang, J., & Long, M. (2021). Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances in Neural Information Processing Systems*, 34.

Xu, Y., & Cohen, S. B. (2018). Stock movement prediction from tweets and historical prices. *Proceedings of ACL*, 1970–1979.

Zeng, A., Chen, M., Zhang, L., & Xu, Q. (2023). Are transformers effective for time series forecasting? *Proceedings of AAAI*, 11121–11128.

Zhang, L., Aggarwal, C., & Kong, X. (2023). A survey on neural network interpretability. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 5(5), 726–742.

Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., & Zhang, W. (2021). Informer: Beyond efficient transformer for long sequence time-series forecasting. *Proceedings of AAAI*, 11106–11115.