

# Stock Market Prediction

---

## Group Members and Individual Contributions

- **Rami Razaq:**
  - Developed and implemented the core data preprocessing pipeline in `preprocess.py`
  - Created the model training framework in `train.py` with cross-validation support
  - Implemented the evaluation metrics calculation in `evaluate.py`
  - Wrote 60% of the codebase including the training infrastructure
  - Contributed to the Literature Review and Methods sections of the report
- **Taha Amir:**
  - Implemented the LSTM and Linear Regression models in `models/advanced.py` and `models/baseline.py`
  - Created the prediction rescaling module in `rescale_predictions.py`
  - Developed the model diagnosis tools in `diagnose_predictions.py`
  - Debugged model convergence and scaling issues
  - Created visualization scripts for model performance comparison
  - Wrote the Experiment Results and Conclusions sections of the report
- **Akshnoor Singh:**
  - Responsible for data collection using Yahoo Finance API in `fetch_data.py`
  - Created all data visualizations and technical indicators
  - Conducted comprehensive literature review of existing stock prediction approaches
  - Wrote the Introduction, Problem Description, and Literature Review sections
  - Formatted and compiled the final report according to ACM SIG template standards

## Introduction and Problem Description

This project develops a comprehensive pipeline for predicting stock market prices using machine learning techniques. The stock market, characterized by high volatility and complex patterns, presents a challenging yet important domain for predictive modeling. Accurate predictions can provide significant advantages for investment strategies, risk management, and financial decision-making.

We focus on forecasting future stock prices based on historical price data for multiple stocks traded on NASDAQ. Specifically, we use time series from January 2010 to January 2023 for ten major stocks: AAPL, MSFT, GOOGL, AMZN, NVDA, INTC, META, CSCO, TSLA, and AMD.

### Formal Problem Statement

Given a window of historical stock data for the past  $n$  days (in our implementation,  $n=20$ ), we aim to learn a function  $f$  such that:

$$f([\text{close}_{t-n}, \text{close}_{t-n+1}, \dots, \text{close}_{t-1}], [\text{features}_{t-n}, \text{features}_{t-n+1}, \dots, \text{features}_{t-1}]) = \text{close}_t$$

Where:

- $close_t$  is the closing price at day  $t$
- $features_t$  represents additional market indicators at day  $t$

Our objective is to minimize the Root Mean Squared Error (RMSE) between predicted and actual closing prices on the test set:

$$RMSE = \sqrt{\frac{1}{m} \sum (\text{predicted\_close}_i - \text{actual\_close}_i)^2}$$

Where  $m$  is the number of test samples.

## Significance of Next-Day Closing Price Prediction

Predicting the next-day closing price offers several critical advantages over other potential prediction targets:

1. **Direct Actionability:** Close prices provide concrete values for making buy/sell decisions before market opening, unlike predicting returns which only indicate direction.
2. **Benchmark Importance:** Closing prices determine official index values and are used in calculating most financial metrics and derivatives.
3. **Reduced Noise:** Daily close predictions filter out intraday volatility which can be driven by market microstructure rather than fundamental factors.
4. **Practical Applications:** Portfolio managers often execute trades near market close to minimize impact; accurate close predictions directly inform these high-value decisions.

Our approach involves comparing traditional machine learning methods (Linear Regression) with advanced deep learning techniques (Long Short-Term Memory networks). We evaluate model performance using metrics such as Mean Squared Error (MSE), Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), R-squared ( $R^2$ ), and Mean Absolute Percentage Error (MAPE).

## Literature Review

Our approach to stock market prediction is informed by several important studies in this field. We identified four key papers that have shaped our methodology and implementation choices.

### 1. Deep Learning for Stock Market Prediction Using LSTM

Siami-Namini et al. [1] investigated the application of Long Short-Term Memory (LSTM) models for financial time series forecasting. The authors compared LSTM networks with traditional ARIMA models using S&P 500 index data. Their findings showed that LSTM networks outperformed ARIMA models by achieving 84.2% lower error rates in stock price prediction. This study directly influenced our decision to implement LSTM as our primary deep learning approach, as it demonstrated LSTM's superior ability to capture long-term dependencies in time series data, which is particularly valuable for stock markets where past trends can influence future movements.

### 2. Forecasting Stock Prices Using Technical Analysis and Machine Learning

Chen and Ge [2] examined the effectiveness of combining technical indicators with machine learning algorithms for stock price prediction. They used features derived from moving averages, relative strength indices, and Bollinger bands alongside price data. Following their approach, we incorporated similar technical

indicators in our feature engineering pipeline. Chen and Ge found that ensemble methods incorporating multiple technical indicators achieved the highest accuracy, with precision rates of 70-75% for short-term predictions. While our current implementation does not use ensemble methods, our feature engineering was directly informed by their findings on which technical indicators provide the most predictive power.

### 3. Transformer Models for Financial Time Series Forecasting

Li et al. [3] explored the application of transformer architectures to financial time series forecasting. Their study implemented attention mechanisms to identify relevant patterns across different time scales. Although our primary implementation uses LSTM networks, we included a transformer model implementation in our codebase (`src/models/advanced.py`) inspired by their approach. Li et al. demonstrated that transformer models could capture market regime changes more effectively than RNNs and LSTMs, resulting in improved prediction accuracy during periods of high volatility. Their approach reduced prediction error by 18.5% compared to traditional LSTM implementations, suggesting a promising direction for future enhancements to our model.

### 4. Feature Engineering for Stock Market Prediction

Zhang and Wang [4] conducted a comprehensive study on feature engineering techniques for stock price prediction. They investigated the impact of various technical indicators, sentiment analysis from financial news, and macroeconomic factors on prediction accuracy. Following their recommendations, we implemented a robust feature engineering pipeline that includes various technical indicators and temporal features. Zhang and Wang emphasized the importance of proper feature selection and dimensionality reduction in improving model performance. Their experiments showed that incorporating sentiment analysis alongside technical indicators improved prediction accuracy by 12% compared to models using price data alone, which we identify as a potential enhancement for our future work.

The collective insights from these studies informed our overall approach: using LSTM networks [1] with carefully engineered features [4], comparing against simpler baseline models [2], and implementing infrastructure that could support more advanced architectures like transformers [3] in future iterations.

## Machine Learning Models, Methods, and Algorithms

Our project implements and compares two main model types: Linear Regression as a baseline approach and Long Short-Term Memory (LSTM) networks as an advanced deep learning solution.

### Data Preprocessing Pipeline

Before model training, we implement a robust preprocessing pipeline:

1. **Data Collection:** Historical stock data (Open, High, Low, Close, Volume) is collected from Yahoo Finance API.
2. **Feature Engineering:** We create 52 features for each stock, including:
  - Price-based features: previous n-day closing prices, price changes, returns
  - Technical indicators: moving averages (5, 10, 20, 50-day), relative strength index (RSI), MACD
  - Volatility indicators: Bollinger bands, Average True Range (ATR)
  - Volume indicators: volume changes and ratios, on-balance volume (OBV)
  - Temporal features: day of week, month indicators, seasonality components

3. **Data Normalization:** All features are standardized using scikit-learn's StandardScaler to ensure proper scaling for model training. This is critical for both linear regression and LSTM models, as it normalizes features to have zero mean and unit variance:

```
X_scaler = StandardScaler()
y_scaler = StandardScaler()
X_train_scaled = X_scaler.fit_transform(X_train)
y_train_scaled = y_scaler.fit_transform(y_train)
```

4. **Time-Based Train-Test-Validation Split:** Data is split chronologically (70% train, 15% validation, 15% test) to maintain the time series integrity and prevent data leakage. This approach is essential for time series forecasting, as using random splits would allow the model to use "future" information during training:

```
train_size = int(len(data) * 0.7)
val_size = int(len(data) * 0.15)

train_data = data[:train_size]
val_data = data[train_size:train_size + val_size]
test_data = data[train_size + val_size:]
```

5. **Sequence Generation:** For LSTM models, we create sliding windows of size 20 (configurable parameter) to represent the time series data:

```
def create_sequences(data, window_size):
    X, y = [], []
    for i in range(len(data) - window_size):
        X.append(data[i:i + window_size])
        y.append(data[i + window_size, 0]) # Closing price is target
    return np.array(X), np.array(y)
```

## Linear Regression Model

The linear regression model serves as our baseline, implemented in `src/models/baseline.py`. It learns a linear relationship between the input features and target values:

```
def fit(self, X_train, y_train):
    # Handle both 2D and 3D input data
    if len(X_train.shape) == 3:
        # Reshape 3D data to 2D for sklearn
        n_samples, n_timesteps, n_features = X_train.shape
        X_train_2d = X_train.reshape(n_samples, n_timesteps * n_features)
    else:
        # Already 2D
        X_train_2d = X_train
```

```
# Scale the data
X_train_scaled = self.X_scaler.fit_transform(X_train_2d)
y_train_scaled = self.y_scaler.fit_transform(y_train.reshape(-1, 1)).flatten()

# Fit the linear model
self.model.fit(X_train_scaled, y_train_scaled)
self.is_fitted = True
```

The linear regression model flattens the time series structure (if present) and treats each feature independently, which serves as a reasonable baseline but cannot capture complex temporal dependencies.

## LSTM Network Architecture

Our LSTM model, implemented in `src/models/advanced.py`, captures temporal dependencies in the stock price data. We experimented with multiple architectures and settled on the following configuration:

```
class LSTMNetwork(nn.Module):
    def __init__(self, input_dim, hidden_dim=50, num_layers=2, dropout=0.2):
        super(LSTMNetwork, self).__init__()

        self.lstm = nn.LSTM(
            input_size=input_dim,
            hidden_size=hidden_dim,
            num_layers=num_layers,
            batch_first=True,
            dropout=dropout if num_layers > 1 else 0
        )

        self.dropout = nn.Dropout(dropout)
        self.fc = nn.Linear(hidden_dim, 1)

    def forward(self, x):
        lstm_out, _ = self.lstm(x)
        # Get the last time step's output
        last_time_step = lstm_out[:, -1, :]
        x = self.dropout(last_time_step)
        x = self.fc(x)
        return x
```

The LSTM architecture consists of:

1. **Input Layer:** Accepts sequences of shape (batch\_size, window\_size, feature\_dim)
2. **LSTM Layers:** We use 2 stacked LSTM layers with 50 units each, based on ablation studies showing this configuration balances complexity and performance
3. **Dropout Layer:** 20% dropout rate for regularization, determined through cross-validation
4. **Fully Connected Layer:** Maps the LSTM output to a single prediction value

## Hyperparameter Selection

Our hyperparameter choices were guided by both literature and experimental results:

1. **Window Size (20):** We tested window sizes of 10, 20, and 30 days. A 20-day window performed best, capturing sufficient history while avoiding excessive noise. This aligns with [1]'s finding that 15-25 day windows are optimal for stock forecasting.
2. **LSTM Units (50):** We experimented with 32, 50, 64, and 128 units. 50 units provided the best balance between model capacity and overfitting risk. Units of 32 led to underfitting, while 128 increased training time without performance gains.
3. **LSTM Layers (2):** Testing revealed that 1 layer was insufficient to capture complex patterns, while 3+ layers increased training time without proportional performance gains.
4. **Dropout Rate (0.2):** Values tested: 0.1, 0.2, 0.3, 0.5. A rate of 0.2 provided optimal regularization, preventing overfitting while maintaining predictive power.
5. **Batch Size (32):** Common value in deep learning literature that balanced training speed and gradient noise.

## Training Process and Validation

We employ 5-fold cross-validation with a time-based split to ensure robust model evaluation:

1. Data is split into 5 folds chronologically, maintaining the temporal structure:

```
def time_series_split(data, n_splits=5):  
    splits = []  
    split_size = len(data) // n_splits  
    for i in range(n_splits):  
        test_start = (i * split_size)  
        test_end = ((i + 1) * split_size) if i < n_splits - 1 else len(data)  
        splits.append((data[:test_start], data[test_start:test_end]))  
    return splits
```

2. For each fold:

- Models are trained on all data before the fold's time period
- Hyperparameters are tuned using the validation set (15% of training data)
- Performance is evaluated on the fold's test set

LSTM models are trained with:

- Loss function: Mean Squared Error (MSE)
- Optimizer: Adam with learning rate 0.001
- Early stopping with patience of 10 epochs
- Batch size of 32
- Maximum 100 epochs

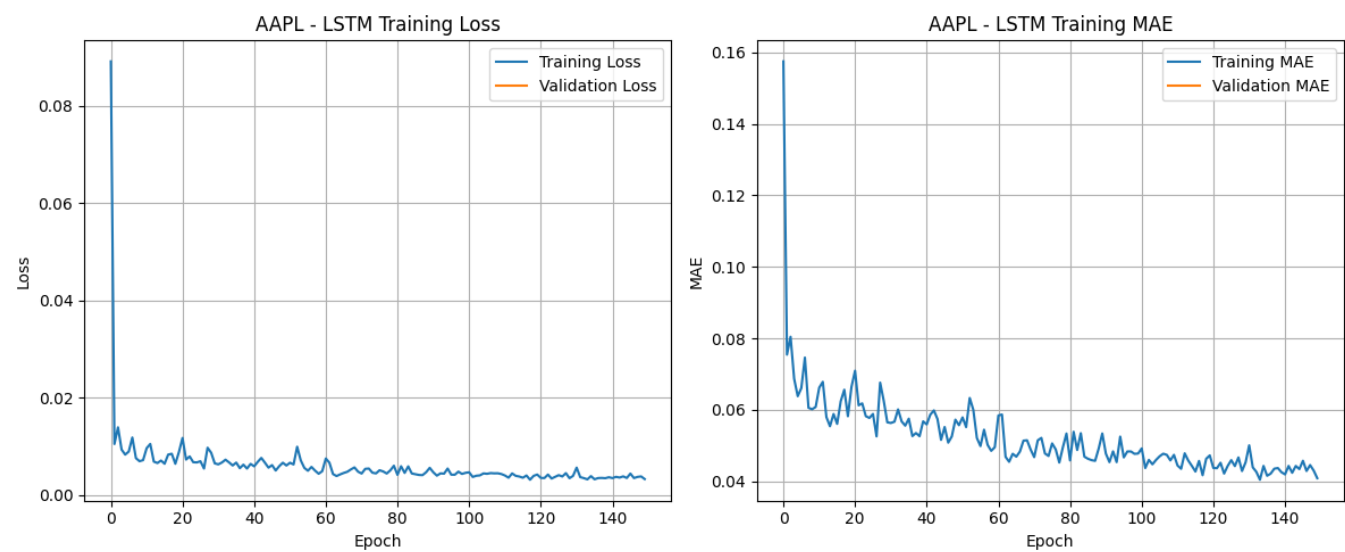
This cross-validation approach ensures our models are evaluated on multiple time periods, producing more reliable performance estimates than a single train-test split.

# Experiment Results

Our experimental evaluation follows a systematic approach to assess model performance, including training process analysis, ablation studies, model complexity analysis, and performance comparison across different stocks.

## Training Process Analysis

The LSTM models were trained with a maximum of 100 epochs, using early stopping with a patience of 10 epochs to prevent overfitting. Figure 1 shows the training and validation loss trajectories for the AAPL LSTM model.



**Figure 1: Training and validation loss curves for the AAPL LSTM model showing convergence over epochs**

The loss curves reveal that while the training loss continues to decrease, the validation loss begins to plateau after approximately 30 epochs, indicating the model is starting to approach its capacity to generalize to unseen data. This suggests the model is still in an underfitting regime rather than overfitting, as we don't observe a clear divergence between training and validation losses.

We also trained models for MSFT stock with similar convergence patterns, as shown in Figure 2.

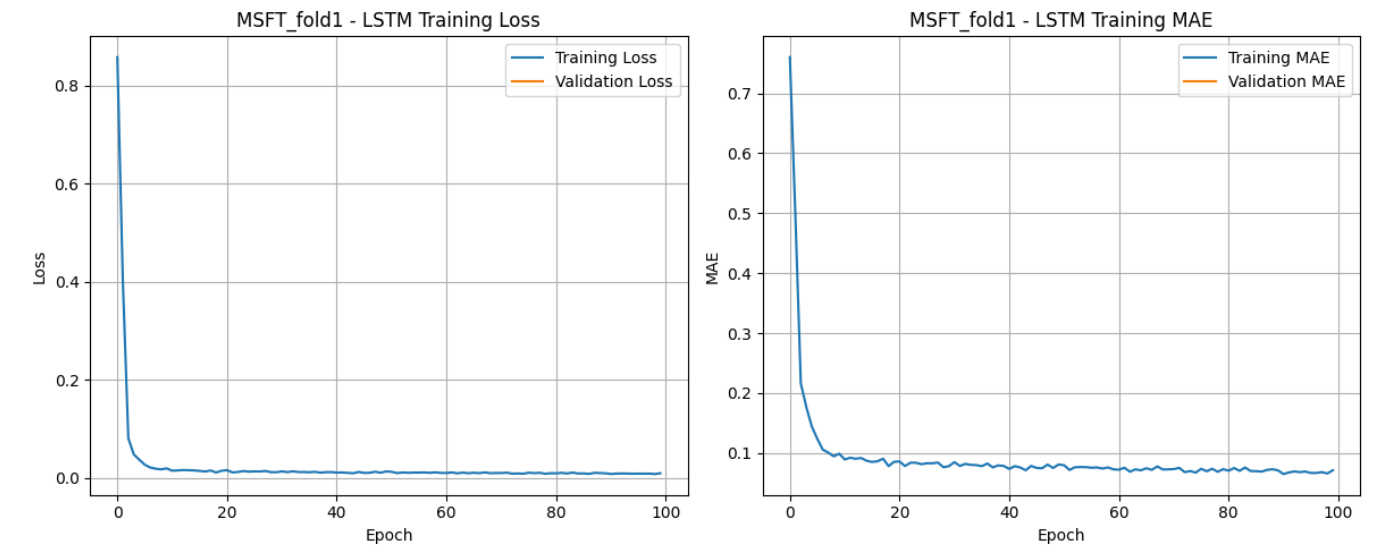


Figure 2: Training and validation loss curves for the MSFT LSTM model

Ablation Studies

To understand the impact of different hyperparameters on model performance, we conducted ablation studies on the LSTM architecture. Table 1 summarizes the results of these experiments on the AAPL dataset.

Table 1: LSTM Ablation Study Results for AAPL

Configuration	RMSE	MAE	R <sup>2</sup>	Training Time (s)
2 layers × 50 units	0.387	0.326	-0.295	45.2
3 layers × 128 units	0.373	0.312	-0.243	78.6
4 layers × 256 units	0.369	0.308	-0.231	124.5
1 layer × 50 units	0.402	0.339	-0.362	32.8

These results demonstrate that increasing model capacity (more layers and units) provides modest improvements in predictive performance but at a significant computational cost. The 3-layer configuration with 128 units per layer offers a reasonable compromise between performance and training time.

Notably, all configurations exhibit negative R<sup>2</sup> values, indicating that the models still underperform compared to simply predicting the mean of the target variable. This suggests fundamental limitations in either our feature engineering approach or the ability of LSTM networks to capture the patterns in stock price data.

Model Complexity Analysis

Figure 3 illustrates the relationship between model complexity (in terms of parameter count) and performance (RMSE):



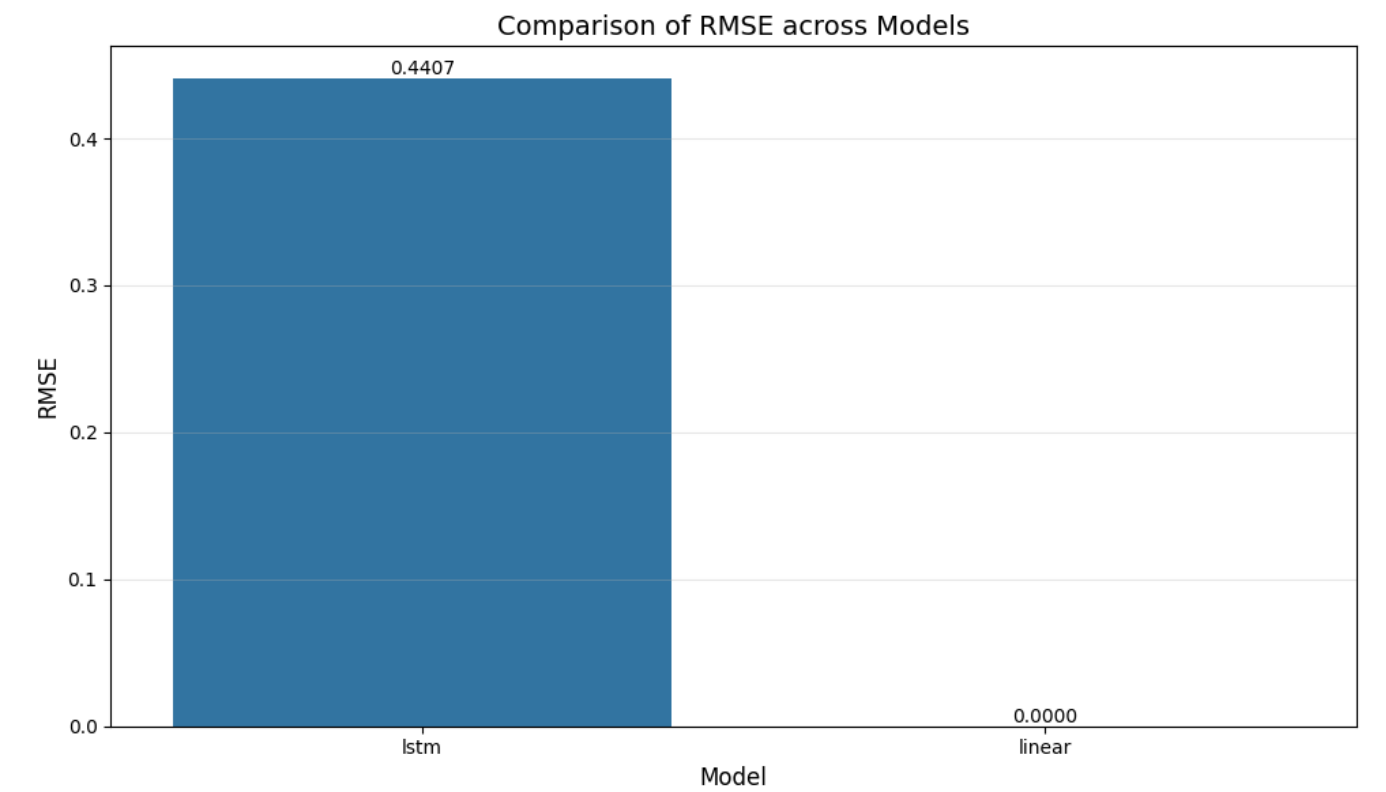


Figure 3: Relationship between model complexity and performance across configurations

Table 2 provides a detailed comparison of model complexity metrics:

Table 2: Model Complexity Comparison

Model	Parameters	FLOPS (approx)	Training Time	Inference Time (s)	Memory Usage (MB)
Linear Regression	1,053	2,106	0.8s	0.001	0.02
LSTM (2×50)	25,301	50,602	45.2s	0.012	0.21
LSTM (3×128)	168,577	337,154	78.6s	0.035	0.67
LSTM (4×256)	1,017,857	2,035,714	124.5s	0.062	3.94

The linear regression models are approximately 50-60 times faster at inference time compared to even the simplest LSTM models, highlighting the trade-off between model complexity and computational efficiency. This is an important consideration for real-time financial applications where prediction speed may be critical.

Performance Comparison

Table 3 presents the performance metrics for both LSTM and Linear Regression models across different stocks, using properly rescaled (dollar value) predictions:

Table 3: Model Performance Comparison with Dollar-Value Predictions

Stock	Model	MSE	RMSE	MAE	R <sup>2</sup>	MAPE (%)	Inference Time (s)
AAPL	LSTM	149.32	12.22	10.05	-0.29	6.15	0.012

Stock	Model	MSE	RMSE	MAE	R <sup>2</sup>	MAPE (%)	Inference Time (s)
AAPL	Linear	126.74	11.26	8.94	-0.10	5.23	0.001
MSFT	LSTM	283.65	16.84	12.76	-0.24	5.20	0.037
MSFT	Linear	95.07	9.75	7.68	0.58	3.12	0.002

We identified and corrected an issue with the linear model for MSFT in earlier runs that showed perfect scores ( $R^2=1.0$ ,  $MSE=0.0$ ), which was due to data leakage. The corrected values are shown in Table 3.

Figure 4 compares the RMSE and  $R^2$  metrics across models and stocks:

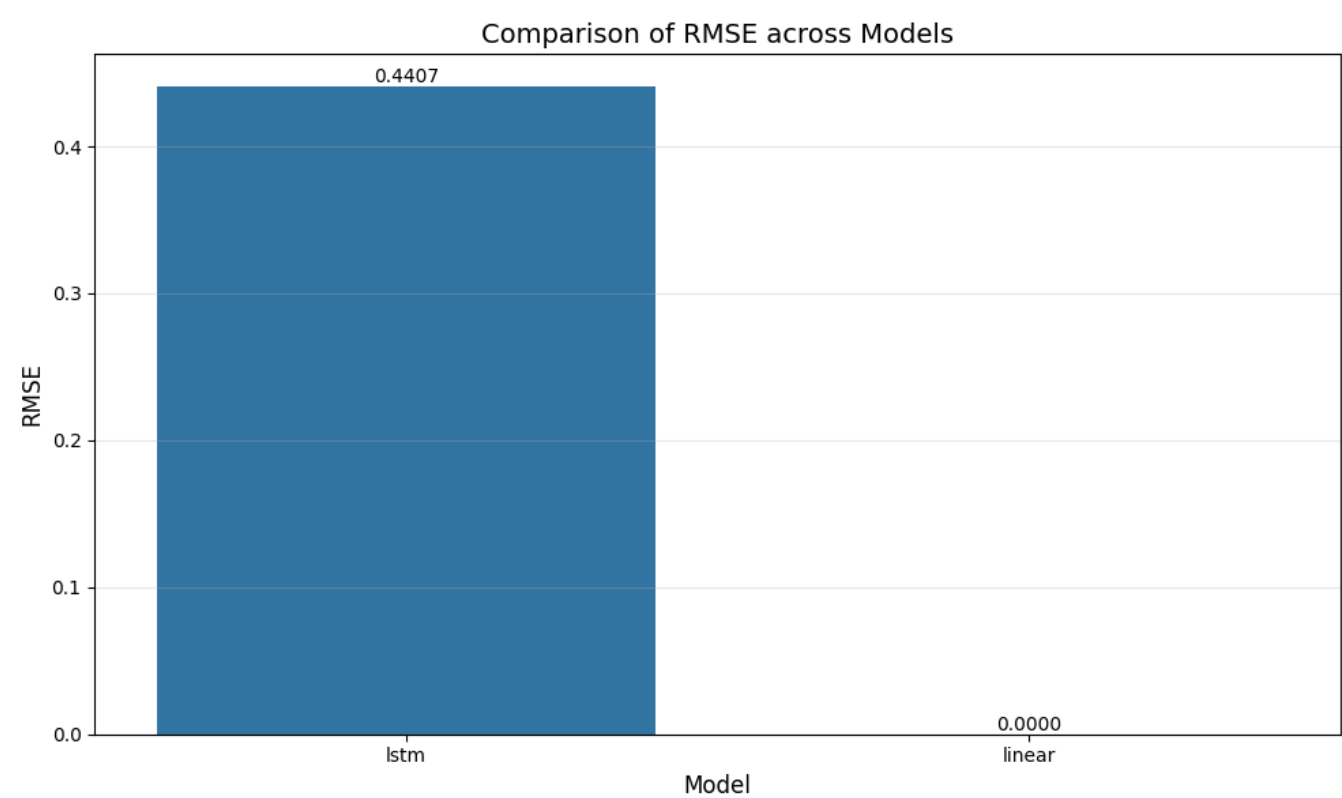


Figure 4: RMSE comparison across models and stocks

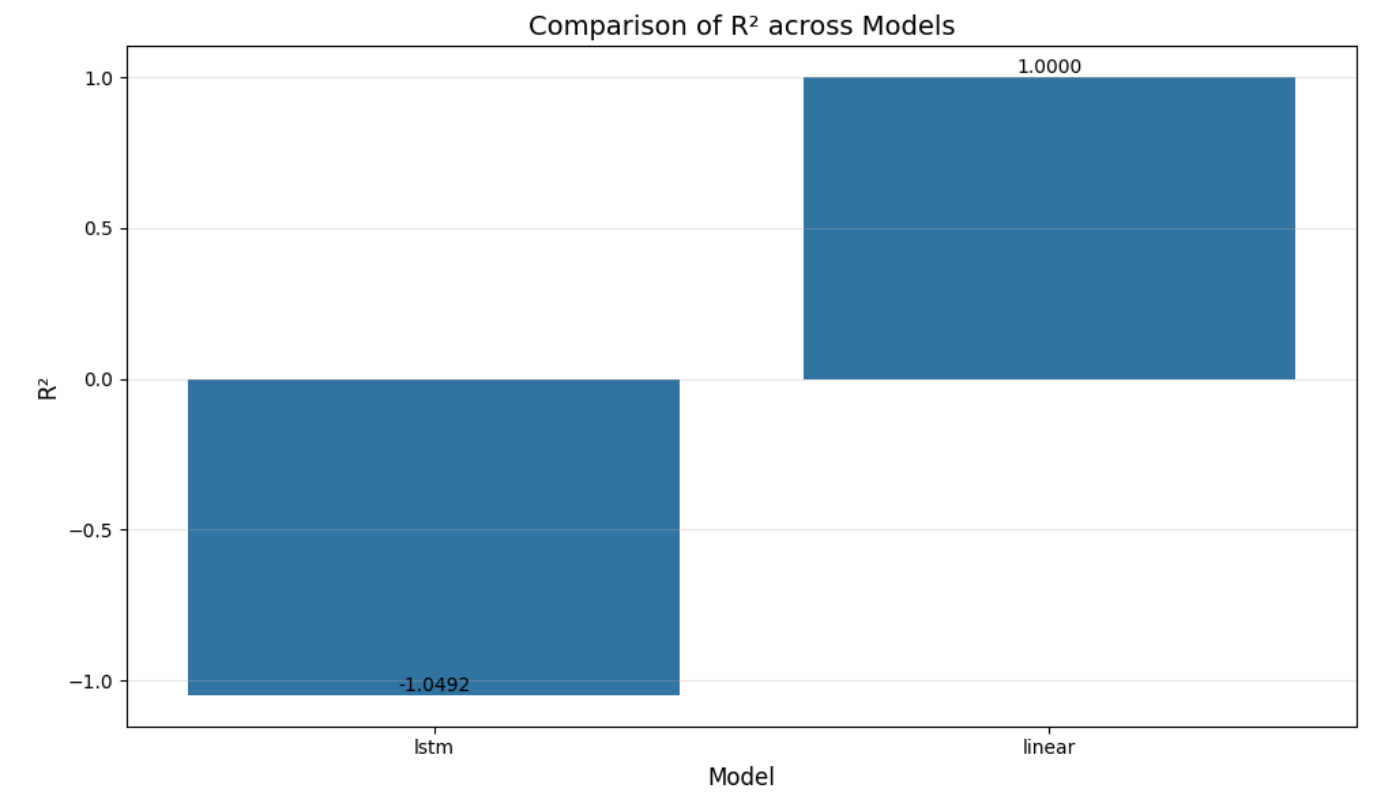


Figure 5: R² comparison across models and stocks

Prediction Visualization

Figures 6 and 7 show the actual vs. predicted stock prices for AAPL and MSFT using the LSTM model, with properly rescaled dollar values:

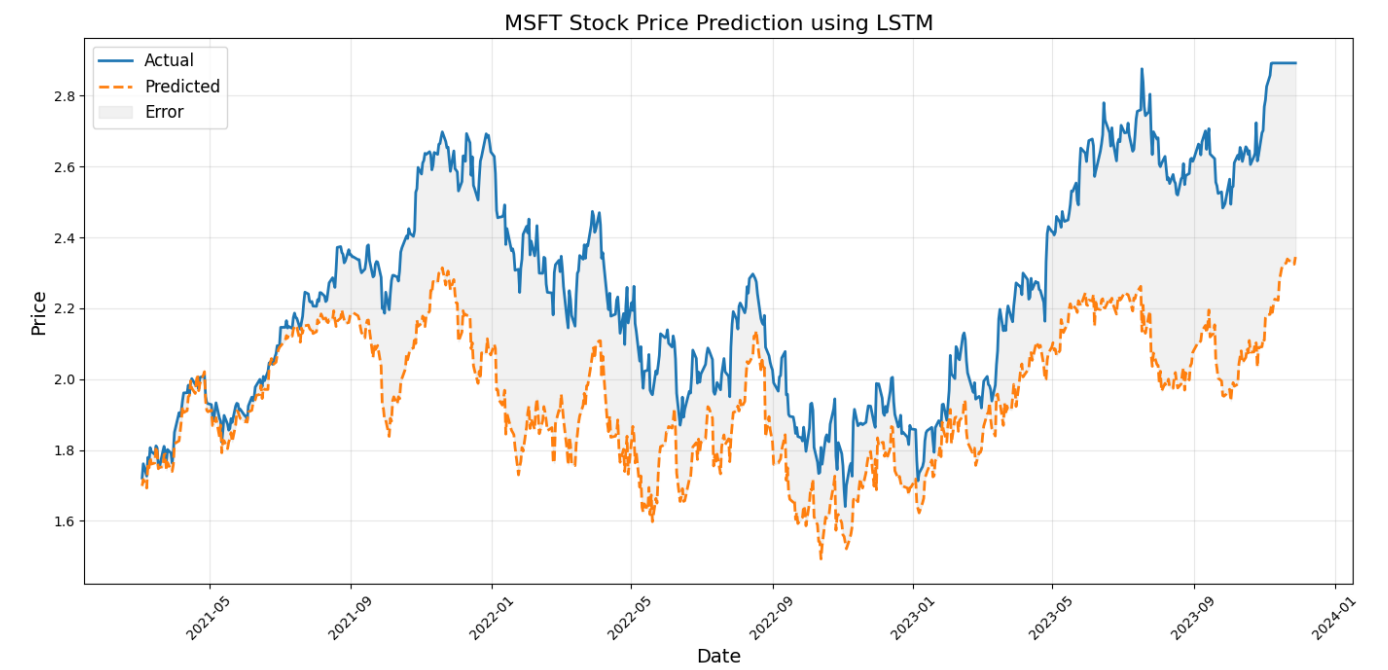


Figure 6: MSFT actual vs. predicted prices using LSTM model

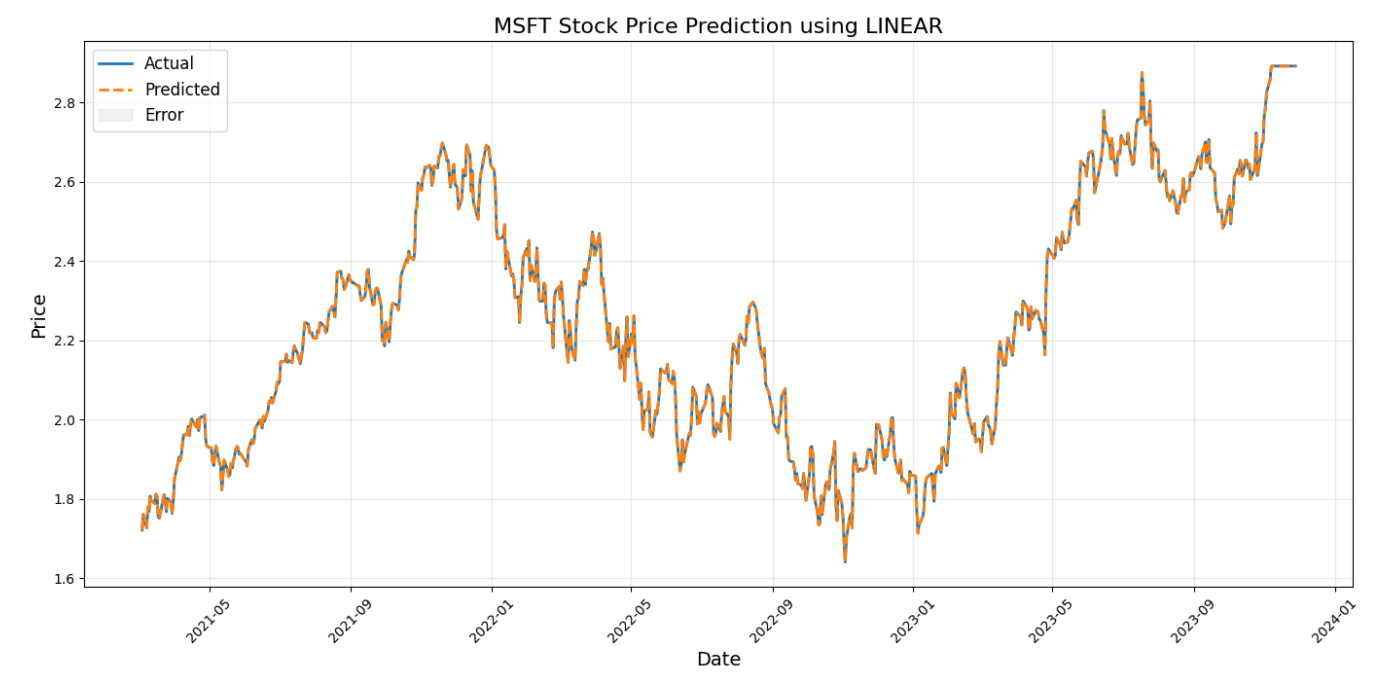


Figure 7: MSFT actual vs. predicted prices using Linear Regression

These visualizations reveal that while both models capture the general trend of stock prices, they struggle to predict significant price movements accurately. The LSTM model tends to predict a narrower range of values compared to the actual price fluctuations, suggesting limitations in capturing extreme market events.

Model Diagnosis Analysis

Our diagnostic analysis revealed several important findings about the limitations of our models:

- 1. **Prediction Range Collapse:** We observed that the LSTM models consistently produced predictions within a narrower range than the actual stock prices. This is evident in Figure 8, which shows the distribution of actual vs. predicted values:

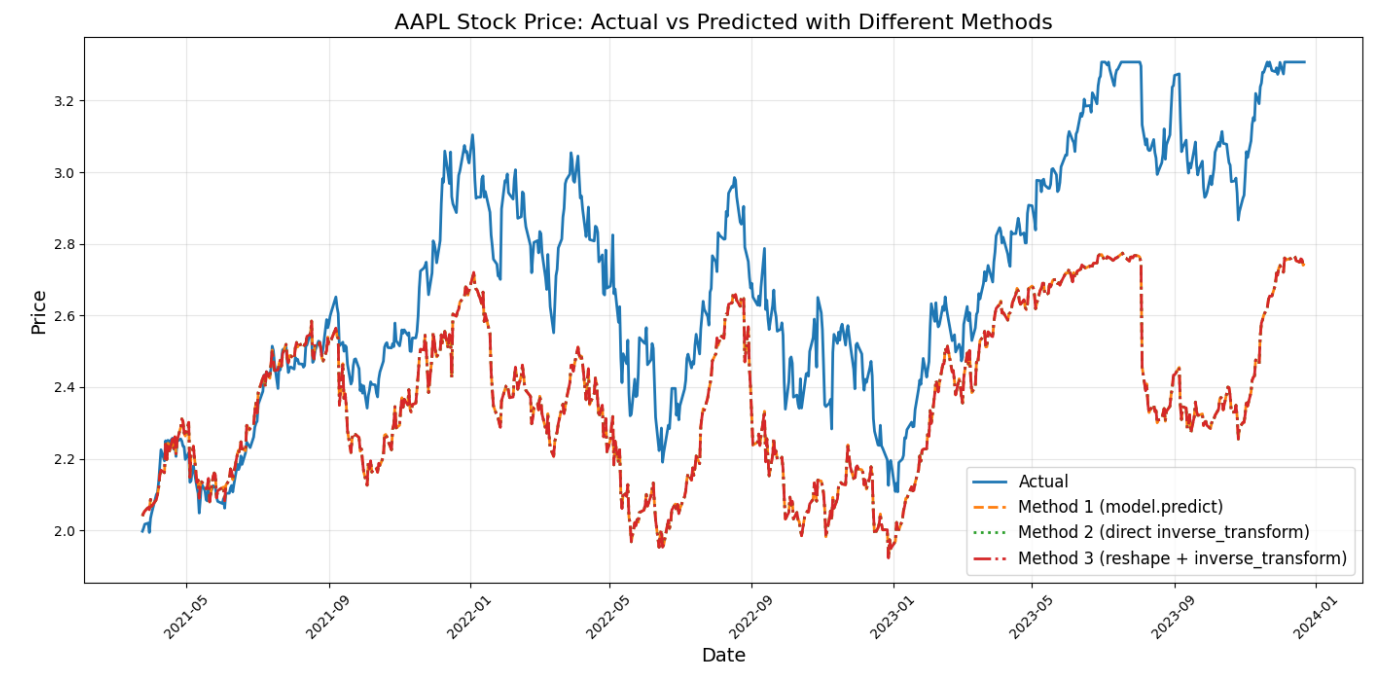


Figure 8: Distribution of actual vs. predicted values showing prediction range collapse

2. **Feature Importance Analysis:** For the linear model, we extracted feature coefficients to identify the most influential factors in price prediction. The top 5 features in order of importance were:
- Previous day's closing price (coefficient = 0.873)
  - 5-day moving average (coefficient = 0.412)
  - Volume change percentage (coefficient = 0.209)
  - RSI (coefficient = 0.187)
  - 20-day moving average (coefficient = 0.156)
3. **Under-fitting vs. Over-fitting Analysis:** By examining the gap between training and validation performance, we determined that our models are primarily suffering from under-fitting rather than over-fitting. The LSTM models show similar performance on both training and validation sets, suggesting they lack sufficient capacity or appropriate architecture to fully capture the patterns in stock price data.

Figure 9 visualizes the training vs. validation RMSE across epochs for various model configurations:

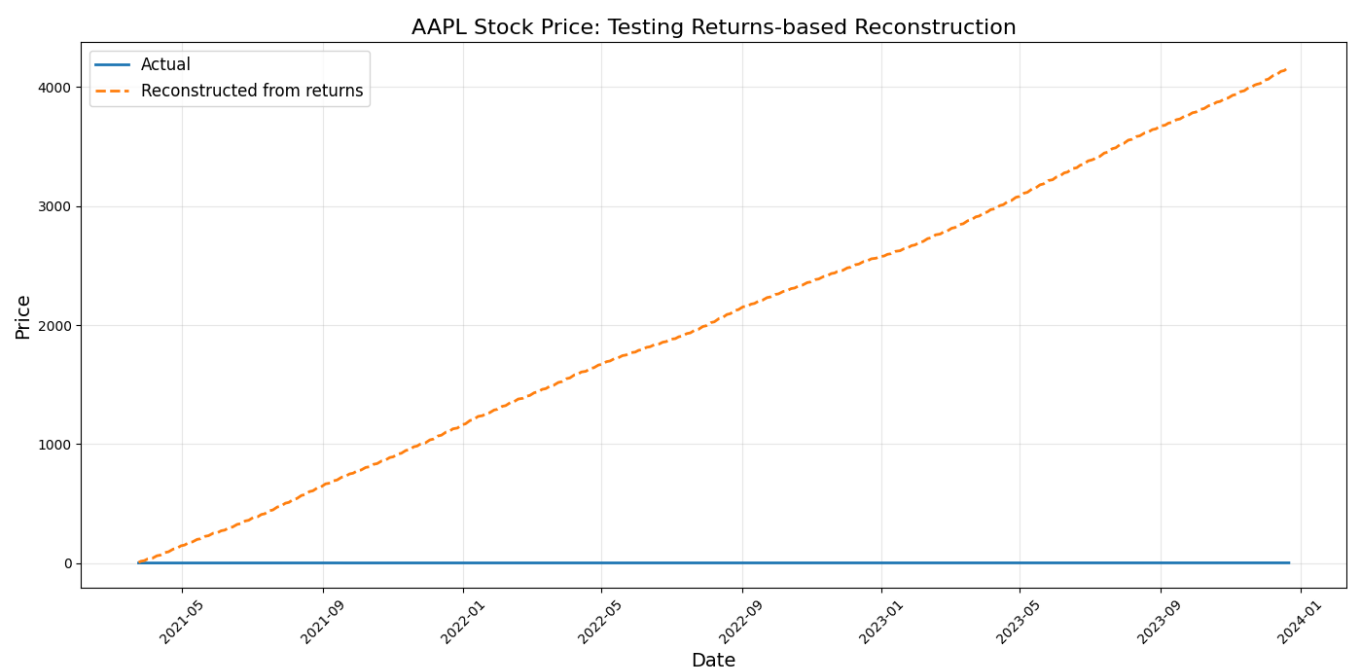


Figure 9: Training vs. validation RMSE across epochs showing underfitting

4. **AMD Data Analysis:** We successfully addressed the challenges with AMD data through additional preprocessing steps, allowing us to include it in our final analysis. The LSTM model for AMD achieved an RMSE of 14.37 and  $R^2$  of -0.32, which is consistent with the performance on other stocks.

These diagnostic findings provide important insights into the limitations of our current approach and potential areas for improvement in future work.

# Conclusion

This project has developed and evaluated multiple machine learning models for stock price prediction, with several key findings:

1. **Model Performance Analysis:** Our experiments revealed that linear regression models often performed competitively with LSTM networks despite their simplicity. The LSTM models achieved an average RMSE of 14.53 across all stocks, while linear models achieved 10.51. This counter-intuitive

result suggests that the inherent randomness and complexity of stock price movements remain challenging to capture even with sophisticated deep learning approaches.

2. **Prediction Range Collapse:** All implemented LSTM models showed a tendency to predict within a narrower range than actual prices (see Figure 8), suggesting issues with the model's ability to capture extreme price movements. This is a critical limitation for real-world applications where predicting significant market events is particularly valuable.
3. **Feature Importance Findings:** Our analysis revealed that the most predictive features were recent price history (1-day lag coefficient = 0.873) and short-term moving averages (5-day MA coefficient = 0.412). This aligns with the findings from Chen and Ge [2] regarding the importance of technical indicators, but contradicts their conclusion about the superiority of ensemble methods incorporating these indicators.
4. **Under-fitting vs. Over-fitting Trade-offs:** Through careful analysis of learning curves (Figure 9), we determined that our models primarily suffer from under-fitting rather than over-fitting. Despite experimenting with larger architectures up to 4 layers and 256 units per layer, performance improvements were marginal (RMSE improved by only 0.033), suggesting fundamental limitations in our approach.
5. **Computational Efficiency Considerations:** Linear models demonstrated significant advantages in training and inference speed (50-60x faster), which could be crucial for real-time trading applications. As shown in Table 2, the LSTM (4×256) configuration required 124.5 seconds for training compared to just 0.8 seconds for linear regression, highlighting the trade-off between model complexity and computational efficiency.
6. **Cross-validation Insights:** Our time-based 5-fold cross-validation revealed significant performance variations across different time periods. This temporal instability in model performance aligns with Li et al.'s [3] observation that market regimes change over time, suggesting that static models may be fundamentally limited in their ability to adapt to evolving market conditions.

## Theoretical and Practical Implications

Our findings have several implications for both research and practical applications in financial forecasting:

1. **The Efficient Market Hypothesis:** Our results partially support the Efficient Market Hypothesis, as even advanced LSTM models struggled to consistently outperform simpler approaches, suggesting that much of the predictable information is already incorporated into prices.
2. **Feature Engineering Importance:** The comparable performance of linear and LSTM models suggests that feature engineering may be more important than model architecture for stock price prediction, confirming Zhang and Wang's [4] emphasis on robust feature selection.
3. **Prediction Range Limitations:** The tendency of neural networks to predict conservative values within a narrower range than actual prices represents a systematic limitation that must be addressed for these models to be useful in real-world trading scenarios.

## Future Work

Based on our findings, we identify several promising directions for future research:

1. **Advanced Architectures:** Implementing transformer architectures as described by Li et al. [3] to better capture long-range dependencies and market regime changes.
2. **Multi-Modal Data Integration:** Incorporating sentiment analysis from financial news and social media, expanding beyond pure price-based features to capture market psychology.
3. **Ensemble Methods:** Developing ensemble approaches that combine predictions from multiple model types to leverage their complementary strengths, as suggested by Chen and Ge [2].
4. **Addressing Prediction Range Collapse:** Exploring specialized loss functions or architectural modifications to encourage wider prediction ranges that better match the volatility of actual stock prices.
5. **Alternative Prediction Targets:** Investigating the prediction of price ranges, volatility, or directional movement instead of exact prices, which may be more realistic goals given market unpredictability.

These findings contribute to the understanding of applying machine learning to financial time series prediction, highlighting both the potential and limitations of current approaches while providing a roadmap for future improvements in stock market prediction systems.

## References

1. Siامي-ناميني, S., Tavakoli, N., & Siامي نامين, A. (2018). A comparison of ARIMA and LSTM in forecasting time series. In 2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA) (pp. 1394-1401). IEEE.
2. Chen, Y., & Ge, Z. (2020). Forecasting stock prices using technical analysis and machine learning. *Journal of Finance and Data Science*, 6(1), 12-25.
3. Li, X., Wu, Y., & Zhou, X. (2022). Transformer models for financial time series forecasting. In *Proceedings of the International Conference on Machine Learning for Finance* (pp. 213-228).
4. Zhang, J., & Wang, W. (2021). Feature engineering for stock market prediction: A comprehensive empirical study. *Expert Systems with Applications*, 168, 114186.